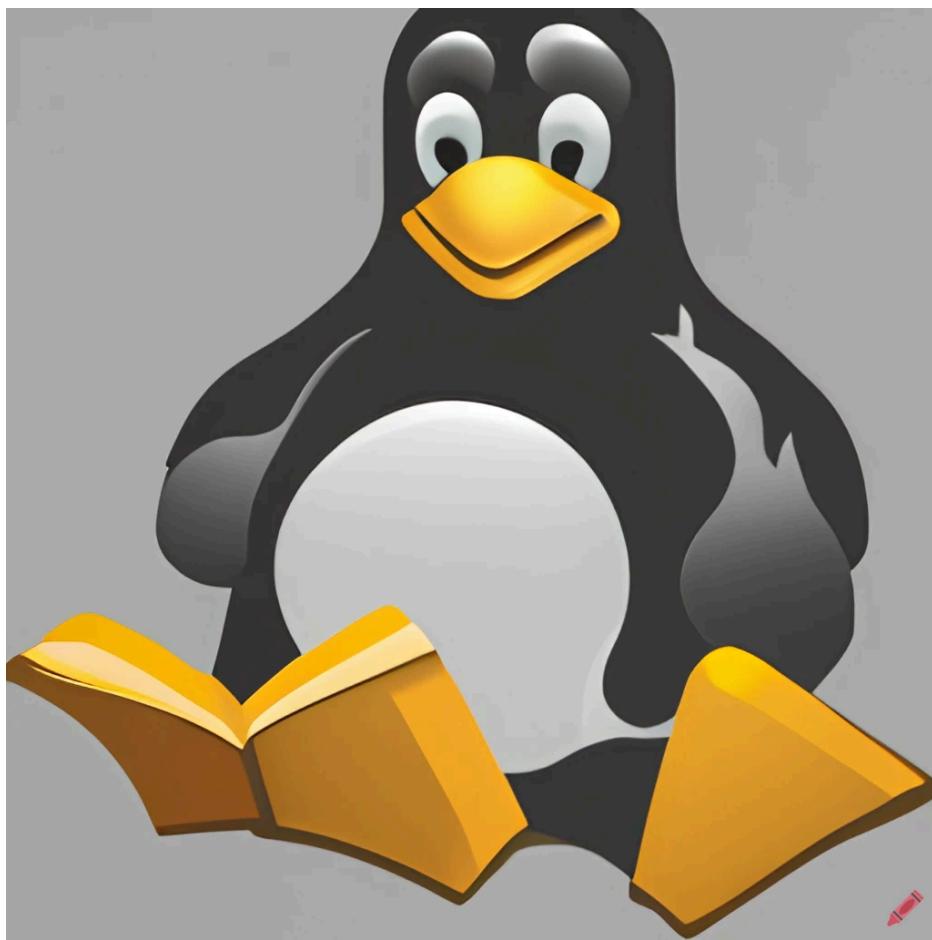


The Linux Process Journey

Version 10.0
October-2025

By Dr. Shlomi Bouthnaru



Created using [Craiyon AI Image Generator](#)

Table of Contents

Table of Contents.....	2
Introduction.....	8
swapper (PID 0).....	9
init (PID 1).....	10
Kernel Threads.....	11
kthreadd (PID 2).....	12
migration.....	14
charger_manager.....	16
idle_inject.....	17
kworker (Kernel Thread Worker).....	18
kdevtmpfs.....	19
cpuhp (CPU Hotplug).....	20
khungtaskd (Kernel Hang Task Daemon).....	20
kswapd.....	22
kcompactd.....	23
md (Multiple Device Driver).....	25
mld (Multicast Listener Discovery).....	27
ksmd (Kernel Same Page Merging).....	28
ttm_swap.....	29
watchdogd (Watchdog Daemon).....	30
zswap-shrink.....	32
khugepaged (Kernel Huge Pages Daemon).....	33
krfcomm (Kernel Radio Frequency Communication Daemon).....	34
ksgxd (Kernel Software Guard eXtensions Daemon).....	35
jbd2 (Journal Block Device 2).....	36
netns (Network Namespace).....	37
oom_reaper (Out-of-Memory Reaper).....	38
kpsmoused (Kernel PS/2 Mouse Daemon).....	39
Slub_flushwq (SLUB Flush Work Queue).....	40
pgdatinit.....	41
kblockd (Kernel Block Daemon).....	42
writeback.....	43
kdiamond (Data Access MONitor).....	44
kintegrityd (Kernel Integrity Daemon).....	45
kthrotld (Kernel Throttling Daemon).....	46
scsi_eh (Small Computer System Interface Error Handling).....	47
blkcg_punt_bio.....	48
napi (New API).....	49

kaudit (Kernel Audit Daemon).....	50
tpm_dev_wq (Trusted Platform Module Device Work Queue).....	51
ipv6_addrconf (IPv6 Address Auto Configuration).....	52
mm_percpu_wq (Per-CPU Memory Work Queue).....	54
inet_frag_wq (IP Fragmentation Work Queue).....	55
kstrp (Stream Parser).....	56
devfreq_wq.....	57
dmcrypt_write (Device Mapper for Transparent Encryption/Decryption).....	58
ModemManager (Modem Management Daemon).....	59
kerneloops.....	60
xargs (Extended Arguments).....	61
cpp (The C Preprocessor).....	62
ntpd (Network Time Protocol Daemon).....	63
gold (The GNU ELF Linker).....	64
strace (System Call Tracing).....	65
ltrace (Library Call Tracer).....	66
kmod (Linux Kernel Module Handling).....	67
lsmod (List Kernel Modules).....	68
insmod (Insert Kernel Module).....	69
rmmmod (Remove Kernel Module).....	70
depmod (Dependency Modules).....	71
cfg80211 (Wireless Configuration).....	72
kdflush (Kernel Device Mapper Flush).....	73
nvme-wq (Non-Volatile Memory Express Work Queue).....	74
] (Checking File Types and Comparing Values).....	75
rcub (Read-Copy Update Boost).....	76
dmesg (Print/Control the Kernel Ring Buffer).....	77
login (Begin Session on The System).....	78
su (Substitute\Switch User).....	79
ps (Process Status).....	80
kacpid (Kernel Advanced Configuration and Power Interface Daemon).....	81
thermald (Thermal Daemon).....	82
dhcpcd (Dynamic Host Configuration Protocol Daemon).....	83
ata_sff (Advanced Technology Attachment Small Form Factor).....	84
uptime (System's Uptime).....	85
agetty (Alternative Linux getty).....	86
Thunar (Xfce File Manager).....	87
ssh (Secure Shell Client).....	88
sshd (OpenSSH SSH Daemon).....	89
passwd (Change User Password).....	90

ssh-agent (Secure Shell Authentication Agent).....	91
hwrng (Hardware Number Random Generator).....	92
psimon (Pressure Stall Information Monitoring).....	93
hv_balloon (Microsoft Hyper-V Memory Ballooning).....	94
snapd (Snap Daemon).....	95
fallocate (Preallocate\Deallocate Space for a File).....	96
ufw (Uncomplicated Firewall).....	97
arpd (Address Resolution Protocol Daemon).....	98
man (Manual Pages).....	99
colord (Color Daemon).....	100
xiccd (X Color Management Daemon).....	101
gnome-terminal.....	102
echo.....	103
lightdm.....	104
wpa_supplicant (Wi-Fi Protected Access client).....	105
wpa_cli (Wi-Fi Protected Access Command Line Client).....	106
wpa_gui (Wi-Fi Protected Access Graphical User Interface).....	107
hostapd (Host Access Point Daemon).....	108
power-profiles-daemon.....	109
mpris-proxy (Media Player Remote Interfacing Specification Proxy).....	110

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	S	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 zeroning? 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 “bpftace” for demonstrating that (if you don’t know about bpftace, 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 bpftace -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 there 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/bus
ybox
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”.

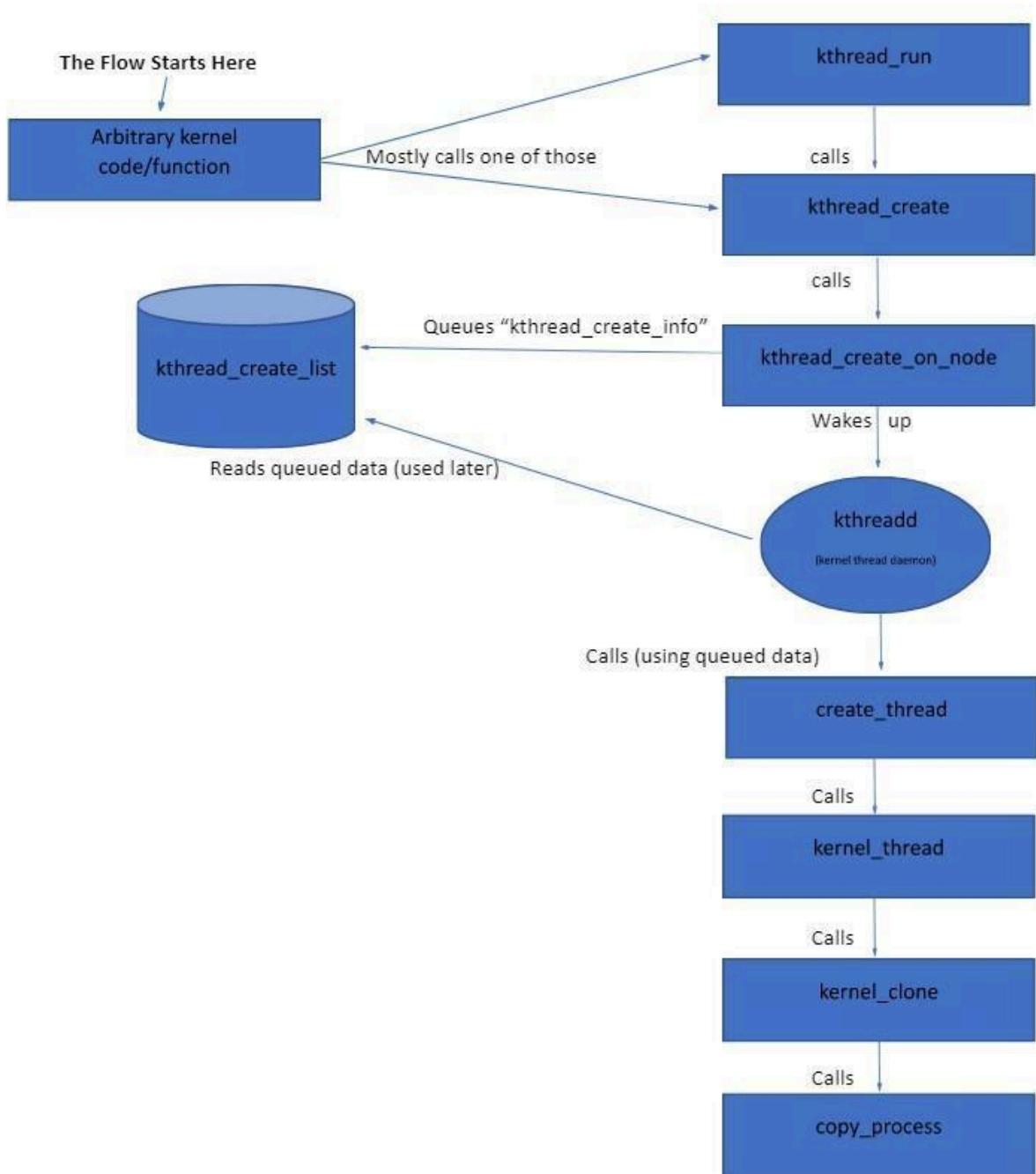
For the sake of simplicity, I have created a flow graph which showcases the flow of creating a kthread, not all the calls are there, only those I thought are important enough. Also, in both cases of macros/functions I used the verb “calls”. The diagram appears at the end of the post. Let me know if it 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-receive 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 udisksd 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”⁹.

⁹ <https://elixir.bootlin.com/linux/latest/source/kernel/sched/core.c#L2325>

charger_manager

The “charger_manager” kernel thread is created by a freezable workqueue¹⁰. Freezable workqueues are basically frozen when the system is moved to a suspend state¹¹. 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 memory¹². The “Charger Manager” kernel module is written by MyungJoo Ham¹³.

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”¹⁴.

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 $
```

¹⁰ <https://elixir.bootlin.com/linux/latest/source/drivers/power/supply/charger-manager.c#L1749>

¹¹ <https://lwn.net/Articles/403891/>

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

¹³ <https://elixir.bootlin.com/linux/latest/source/drivers/power/supply/charger-manager.c#L1768>

¹⁴ <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 kwoker 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 kernel¹⁵). 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_ldisc
kworker/u2:2-46624 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
```

¹⁵ <https://www.kernel.org/doc/Documentation/trace/ftrace.txt>

kdevtmpfs

“kdevtmpfs” is a kernel thread which was created using the “kthread_run” function¹⁶. “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 filesystem¹⁷.

We can see the specific line of code that is used in order to create the mounting point “/dev”¹⁸. The mountpoint is created using the function “init_mount”¹⁹. 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 code²⁰.

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 noram
```

¹⁶ <https://elixir.bootlin.com/linux/v6.2-rc1/source/drivers/base/devtmpfs.c#L474>

¹⁷ <https://elixir.bootlin.com/linux/v6.2-rc1/source/drivers/base/devtmpfs.c#L3>

¹⁸ <https://elixir.bootlin.com/linux/v6.2-rc1/source/drivers/base/devtmpfs.c#L377>

¹⁹ <https://elixir.bootlin.com/linux/v6.2-rc1/source/fs/init.c#L16>

²⁰ <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"²¹.

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"²² is called. In order to bring up a CPU function "cpu_device_up"²³ 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
```

²¹ <https://elixir.bootlin.com/linux/latest/source/kernel/cpu.c#L65>

²² <https://elixir.bootlin.com/linux/latest/source/kernel/cpu.c#L1225>

²³ <https://elixir.bootlin.com/linux/latest/source/kernel/cpu.c#L1439>

khungtaskd (Kernel Hang Task Daemon)

This kernel thread “`khungtaskd`” is used in order to help with identifying and debugging “Hung Tasks”. This kernel thread is scheduled every 120 seconds (that is the default value). We can say “`khungtaskd`” is used for detecting tasks which are stuck in uninterruptible sleep (state “D” in ps output). We can also go over the code of the kernel thread as part of the Linux kernel source code²⁴.

The basic algorithm of “khungtaskd” is as follows: Iterate over all running tasks on the system and if there are ones marked as TASK_UNINTERRUPTIBLE and it was scheduled at least once in the last 120 seconds it is considered as hung. When a task is considered hung it’s “call stack” is dumped and if the CONFIG_LOCKDEP is also enabled then all of the locks held by the tasks are outputed also.

If we want we can change the sampling interval using the sysctl interface, “/proc/sys/kernel/hung_task_timeout_secs”. We can also verify that the default is 120 seconds by reading it - as shown in the screenshot below.

In order to demonstrate the operation of “khungtaskd” I have executed the following bpftrace one liner - “`sudo bpftrace -e 'kfunc:check_hung_uninterruptible_tasks { printf("%s:%d\n",curtask->comm,curtask->pid); }'`”. The trace prints the name of the task and it’s pid when the function “`check_hung_uninterruptible_tasks`” is called²⁵ - You can see the output in the screenshot below.

²⁴ https://elixir.bootlin.com/linux/latest/source/kernel/hung_task.c

https://elixir.bootlin.com/linux/latest/source/kernel/hung_task.c#L178

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”²⁶. The code of the page allocator is located at “mm/page_alloc.c”²⁷.

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()”²⁸. 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).

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

²⁷ https://elixir.bootlin.com/linux/latest/source/mm/page_alloc.c

²⁸ <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 born²⁹.

“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 migration³⁰ to do all the heavy lifting³¹. In order for “kcompactd” to work we should compile the kernel with “CONFIG_COMPACTION” enabled. Also, when a Linux system identifies that it is tight low in available memory the “kcompactd” won’t perform memory compaction memory³².

Overall, the “kcompactd” kernel thread is created in “kcompactd_run” function³³ which is called by “kcompactd_init”³⁴. The function “kcompactd_init” is started by “subsys_initcall”³⁵, which is responsible for initializing a subsystem.

The kernel thread starts the function “static int kcompactd(void *p)”³⁶. An instance of the kernel thread is created for each node (like vCPU) on the system³⁷. 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 kswapd³⁸. 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”³⁹ 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.

²⁹ <https://lwn.net/Articles/368869/>

³⁰ <https://lwn.net/Articles/157066/>

³¹ <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L5>

³² <https://www.linux-magazine.com/Issues/2015/179/Kernel-News>

³³ <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L2996>

³⁴ <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L3048>

³⁵ <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L3065>

³⁶ <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L2921>

³⁷ <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L3061>

³⁸ <https://www.slideshare.net/AdrianHuang/memory-compaction-in-linux-kernelpdf>

³⁹ <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 workqueue⁴⁰. 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 filesystem⁴¹.

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)⁴².

Overall, “md” supports different RAID types: RAID 1 (mirroring), RAID 4, RAID 5, RAID 6 and RAID 10⁴³. Besides that, “md” also supports pseudo RAID technologies like: RAID 0, LINAR, MULTIPATH and FAULTY⁴⁴.

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 $
```

⁴⁰ <https://elixir.bootlin.com/linux/v6.1/source/drivers/md/md.c#L9615>

⁴¹ <https://linux.die.net/man/8/mdadm>

⁴² <https://man7.org/linux/man-pages/man4/md.4.html>

⁴³ <https://www.prepressure.com/library/technology/raid>

⁴⁴ 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])⁴⁵. 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”⁴⁶. There is also the command line utility “mdadm” that can help with managing those devices⁴⁷.

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”)⁴⁸. More information about initcalls will be included on a future writeup.

⁴⁵ <https://www.oreilly.com/library/view/managing-raid-on/9780596802035/ch01s03.html>

⁴⁶ <https://raid.wiki.kernel.org/index.php/Mdstat>

⁴⁷ <https://linux.die.net/man/8/mdadm>

⁴⁸ <https://elixir.bootlin.com/linux/v6.1/source/drivers/md/md.c#L9947>

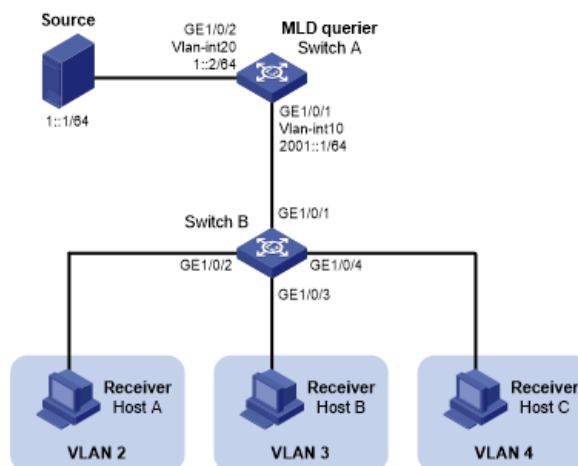
mld (Multicast Listener Discovery)

“mld” is a kernel thread which was created using a workqueue⁴⁹. 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 Windows⁵⁰ and Linux⁵¹.

We can think about it like IGMP⁵² 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 types⁵³.

Overall, MLD has three major message types: “Multicast Listener Query”, “Multicast Listener Report” and “Multicast Done”⁵⁴. For more information about them I suggest reading the following link. Also, a more detailed explanation about the different MLD operations can be found in <https://ipcisco.com/lesson/mld-operations/>.

What “mld” does is to send MLD report messages⁵⁵ which are sent by an MLD host (see the diagram below⁵⁶) and processes messages⁵⁷. From the source code we can see that there are definitions for structs representing both MLDv1 and MLDv2 headers.



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

⁵⁰ <https://learn.microsoft.com/en-us/windows/win32/winsock/igmp-and-windows-sockets>

⁵¹ <https://lwn.net/Articles/29489/>

⁵² <https://www.cloudflare.com/learning/network-layer/what-is-igmp/>

⁵³ <https://www.ibm.com/docs/en/zos/2.2.0?topic=protocol-multicast-listener-discovery>

⁵⁴ <https://community.cisco.com/t5/networking-knowledge-base/multicast-listener-discovery-mld/ta-p/3112082>

⁵⁵ <https://elixir.bootlin.com/linux/latest/source/net/ipv6/mcast.c#L3185>

⁵⁶ https://techhub.hpe.com/eginfolib/networking/docs/switches/5130ei/5200-3944_ip-multi_cg/content/images/image33.png

⁵⁷ <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(). We can also go over the source code of “ksmd” source code as part of the Linux kernel⁵⁸.

The way “ksmd” works is as follows. Scanning main memory for frames (“physical pages”) holding identical data and collects 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_run⁵⁹. 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
max_page_sharing    pages_sharing   pages_volatile  stable_node_chains
merge_across_nodes  pages_to_scan   run             stable_node_chains_prune_millisecs
                                         stable_node_dups
                                         use_zero_pages
```

⁵⁸ <https://elixir.bootlin.com/linux/v6.13.7/source/mm/ksm.c>

⁵⁹ <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 objects⁶⁰.

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)⁶¹.

Moreover, “ttm_swap” is a single threaded workqueue as seen in the Linux source code⁶². 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 driver⁶³.

TTM is at the end a kernel module, you can find the source code and the Makefile in the kernel source tree⁶⁴. Based on the module source code it is written by Thomas Hellstrom and Jerome Glisse⁶⁵. Also, it is described as “TTM memory manager subsystem (for DRM device)”⁶⁶. 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
```

⁶⁰ <https://docs.kernel.org/gpu/drm-mm.html>

⁶¹ <https://docs.kernel.org/gpu/drm-internals.html>

⁶² https://elixir.bootlin.com/linux/v5.12.19/source/drivers/gpu/drm/ttm/ttm_memory.c#L424

⁶³ <https://www.systutorials.com/docs/linux/man/7-drm-ttm/>

⁶⁴ <https://elixir.bootlin.com/linux/v6.1-rc2/source/drivers/gpu/drm/ttm>

⁶⁵ https://elixir.bootlin.com/linux/v6.1-rc2/source/drivers/gpu/drm/ttm/ttm_module.c#L89

⁶⁶ https://elixir.bootlin.com/linux/v6.1-rc2/source/drivers/gpu/drm/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 support⁶⁷.

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 restart⁶⁸.

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 timer⁶⁹. When the timer gets to a specific value (configured ahead) and it is not elapsed by the software the hardware will restart the system⁷⁰.

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 hardware⁷¹. Examples for different watchdog drives/modules for specific chips can be found in the source tree of Linux kernel⁷². 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. 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).

⁶⁷ <https://github.com/torvalds/linux/blob/master/Documentation/watchdog/watchdog-api.rst>

⁶⁸ <https://linuxhint.com/linux-kernel-watchdog-explained/>

⁶⁹ https://en.wikipedia.org/wiki/Watchdog_timer

⁷⁰ <https://developer.toradex.com/linux-bsp/how-to/linux-features/watchdog-linux/>

⁷¹ <https://linux.die.net/man/8/watchdog>

⁷² <https://elixir.bootlin.com/linux/v6.0.11/source/drivers/watchdog>

Lastly, if you want to see the status of the watchdog you can use the command “wdctl”⁷³ - As can be seen in the screenshot below⁷⁴.

```
[root@ako-kaede-mirai]-(12:25am--09/06) ~-~  
[root@kousekip] ~ 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
```

⁷³ <https://man7.org/linux/man-pages/man8/wdctl.8.html>

⁷⁴ 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\swapfile⁷⁵. 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 pool⁷⁶.

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 cache⁷⁷. The “zswap_frontswap_store” is the function that attempts to compress and store a single page⁷⁸.

The kernel thread “zswap-shrink” is created based on a workqueue⁷⁹. 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 -efl 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 zswap⁸⁰. We can also read the parameters on “/sys/module/zswap/parameters”.

⁷⁵ <https://www.kernel.org/doc/html/v4.18/vm/frontswap.html>
⁷⁶ <https://elixir.bootlin.com/linux/latest/source/mm/zswap.c>
⁷⁷ <https://www.kernel.org/doc/html/v4.18/vm/zswap.html>
⁷⁸ <https://elixir.bootlin.com/linux/v6.1-rc2/source/mm/zswap.c#L1097>
⁷⁹ <https://elixir.bootlin.com/linux/v6.1-rc2/source/mm/zswap.c#L1511>
⁸⁰ <https://wiki.archlinux.org/title/zswap>

khugepaged (Kernel Huge Pages Daemon)

The kernel thread “kugepaged” is created using the “kthread_run()” function⁸¹. It is responsible for the “Transparent Hugepage Support” (aka THP). “kugepaged” scans memory and collapses sequences of basic pages into huge pages⁸².

We can manage and configure TPH using sysfs⁸³ or by using the syscalls “madvise”⁸⁴ and “prctl”⁸⁵. The scan of memory is done by calling “khugepaged_do_scan()”⁸⁶ which in turn calls “khugepaged_scan_mm_slot()”⁸⁷. 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 mapping⁸⁸.

```
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
```

⁸¹ <https://elixir.bootlin.com/linux/latest/source/mm/khugepaged.c#L2551>
⁸² <https://www.kernel.org/doc/html/latest/admin-guide/mm/transhuge.html>
⁸³ <https://www.kernel.org/doc/html/latest/admin-guide/mm/transhuge.html#thp-sysfs>
⁸⁴ <https://man7.org/linux/man-pages/man2/madvise.2.html>
⁸⁵ <https://man7.org/linux/man-pages/man2/prctl.2.html>
⁸⁶ <https://elixir.bootlin.com/linux/latest/source/mm/khugepaged.c#L2404>
⁸⁷ <https://elixir.bootlin.com/linux/v6.1.12/source/mm/khugepaged.c#L2250>
⁸⁸ <https://www.kernel.org/doc/html/latest/admin-guide/mm/transhuge.html>

krfcomm (Kernel Radio Frequency Communication Daemon)

“krfcomm” is a kernel which is started by executing “kthread_run()” function⁸⁹. The kernel thread executes the “rfcomm_run()” function⁹⁰. Thus, we can say that “krfcomm” is responsible for RFCOMM connections⁹¹.

RFCOMM (Radio Frequency Communication)⁹² 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 Bluetooth⁹³.

Moreover, there is also an “rfcomm” cli tool in Linux. It is used to inspect and maintain RFCOMM configuration⁹⁴. We can also go over the protocol specification⁹⁵. 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 them⁹⁶.

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:   2787EEC4EC282A1A24A7701
depends:      bluetooth
retpoline:    y
intree:       y
name:        rfcomm
vernmagic:   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:B8: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:E2:8B:98:E9:4A:50:
              77:00:37:08:0A:12:CD:5D:84:F7:2F:4A:FA:CB:58:6B:B9:C4:7B:C0:
              08:1C:EC:61:33:FA:7E:A6:69:6B:FD:E7:02:30:69:C8:06:98:12:9C:
              E3:B3:25:33:03:12:81:D6:77:S9:54:F5:8E:5B:D5:FF:C4:5D:D1:FI:
              02:0E:16:6B:2E:33:84:97:2D: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)
```

⁸⁹ <https://elixir.bootlin.com/linux/latest/source/net/bluetooth/rfcomm/core.c#L2215>

⁹⁰ <https://elixir.bootlin.com/linux/latest/source/net/bluetooth/rfcomm/core.c#L2109>

⁹¹ <https://stackoverflow.com/questions/57152408/what-is-the-internal-mechanics-of-socket-function>

⁹² <https://www.bluetooth.org/rfcomm.htm>

⁹³ https://en.wikipedia.org/wiki/List_of_Bluetooth_protocols

⁹⁴ <https://linux.die.net/man/1/rfcomm>

⁹⁵ <https://www.bluetooth.com/specifications/specs/rfcomm-1-1/>

⁹⁶ https://www.amd-e-technik.uni-rostock.de/ma/gol/lectures/wirlec/bluetooth_info/rfcomm.html

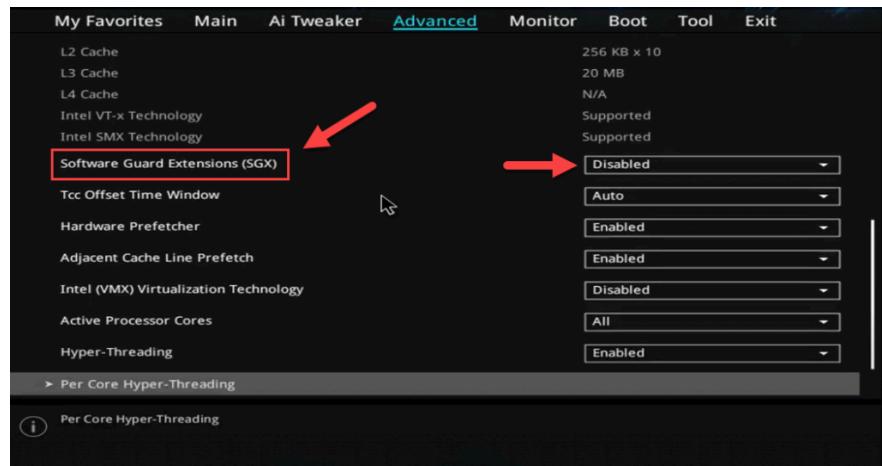
ksgxd (Kernel Software Guard eXtensions Daemon)

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 regions⁹⁷. For more information about SGX you can read my writeup about it⁹⁸.

“ksgxd” is a kernel which is started by executing “kthread_run()” function⁹⁹. The kernel thread executes the “ksgxd” function¹⁰⁰. “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 memory¹⁰¹ like “kswapd”¹⁰².

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 below¹⁰³.

Lastly, there is a great guide for an example SGX app using a Linux VM on Azure that I encourage you to read¹⁰⁴. There is also specific documentation\information regarding the Linux stack for SGX¹⁰⁵. By the way, we can go over the following Intel’s GitHub repository dedicated to SGX¹⁰⁶.



⁹⁷ <https://docs.kernel.org/x86/sgx.html>

⁹⁸ <https://medium.com/@boutnaru/security-sgx-software-guard-extension-695cab7dbcb2>

⁹⁹ <https://elixir.bootlin.com/linux/v6.1.10/source/arch/x86/kernel/cpu/sgx/main.c#L427>

¹⁰⁰ <https://elixir.bootlin.com/linux/v6.1.10/source/arch/x86/kernel/cpu/sgx/main.c#L395>

¹⁰¹ <https://elixir.bootlin.com/linux/v6.1.10/source/arch/x86/kernel/cpu/sgx/main.c#L188>

¹⁰² <https://medium.com/@boutnaru/the-linux-process-journey-kswapd-22754e783901>

¹⁰³ <https://phoenixnap.com/kb/intel-sgx>

¹⁰⁴ <https://tsmatz.wordpress.com/2022/05/17/confidential-computing-intel-sgx-enclave-getting-started/>

¹⁰⁵ https://download.01.org/intelsgxstack/2021-12-08/Getting_Started.pdf

¹⁰⁶ <https://github.com/intel/linux-sgx>

jbd2 (Journal Block Device 2)

“JBD” stands for “Journal Block Device”¹⁰⁷. “jbd2” is a kernel which is started by executing “kthread_run()” function¹⁰⁸. 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 log¹⁰⁹.

The kernel thread executes the “kjournald2()” function¹¹⁰. 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 file¹¹¹.

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 ext4¹¹².

```
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:
               ?D:7B:59:12:58:D2:C5:5D:99:63:81:6B:60:D2:63:6C:0F:18:5A:26:
               9D:93
```

¹⁰⁷ <https://manpages.ubuntu.com/manpages/jammy/man1/pmdajbd2.1.html>

¹⁰⁸ <https://elixir.bootlin.com/linux/v6.2.1/source/fs/jbd2/journal.c#L277>

¹⁰⁹ <https://elixir.bootlin.com/linux/v6.2.1/source/fs/jbd2/journal.c#L169>

¹¹⁰ <https://elixir.bootlin.com/linux/v6.2.1/source/fs/jbd2/journal.c#L152>

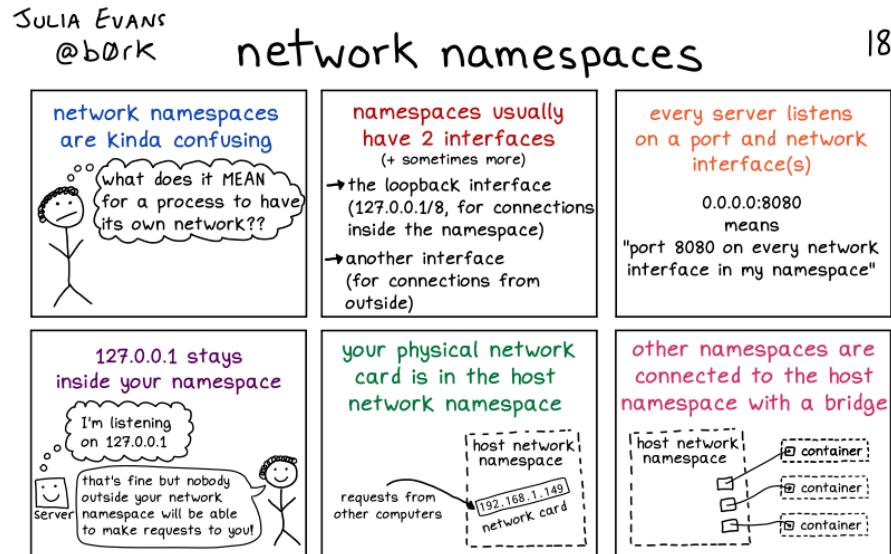
¹¹¹ https://en.wikipedia.org/wiki/Journaling_block_device

netns (Network Namespace)

The kernel thread “netns” is based on a single threaded workqueue¹¹³, which is created when the network namespace¹¹⁴ is initialized (net_ns_init()). Also, for a reminder you can also check out the diagram below¹¹⁵.

“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()” function¹¹⁶.

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()” function¹¹⁷.



¹¹³ https://elixir.bootlin.com/linux/v6.2-rc4/source/net/core/net_namespace.c#L1106

¹¹⁴ <https://medium.com/@boutnaru/linux-namespaces-network-namespace-part-3-7f8f8e06fef3>

¹¹⁵ <https://wizardzines.com/comics/network-namespaces/>

¹¹⁶ https://elixir.bootlin.com/linux/v6.2.3/source/net/core/net_namespace.c#L565

¹¹⁷ 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” function¹¹⁸. Basically, it is the implementation of the OMM (Out-of-Memory) killer¹¹⁹ function of the Linux kernel. The function which is executed by the thread is “oom_reaper”¹²⁰ which calls “oom_reap_task”¹²¹.

Overall, based on the documentation the goal of the “oom_reaper” kernel thread is to try and reap the memory used by the OOM victim¹²². “oom_reaper” sleeps until it is waked up¹²³ which is after OOM kills the process¹²⁴.

Lastyl, after killing the process the victim is queued so the “oom_reaper” can release the resources¹²⁵. You can see an example of the log created by OOM after killing a process¹²⁶.

```
4i7kwlagZ kernel: [ 1440]      0 1440    24585      250      0
4i7kwlagZ kernel: [ 1442]      0 1442    27085      100      0
4i7kwlagZ kernel: [ 1458]      0 1458    85095     8955      0
4i7kwlagZ kernel: Out of memory: Kill process 3223 (java) score
4i7kwlagZ kernel: Killed process 3223, UID 0, (java) total-vm:1
4i7kwlagZ yum[1458]: Updated: 7:squid-3.1.23-24.el6.x86_64
4i7kwlagZ squid[1511]: Squid Parent: child process 1513 started
```

¹¹⁸ https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L735

¹¹⁹ <https://medium.com/@bouthnari/linux-out-of-memory-killer-oom-killer-bb2523da15fc>

¹²⁰ https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L640

121 https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L609122 https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L504123 https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L680124 https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L947125 https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L992126 <https://blog.cadata.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 workqueue¹²⁷ 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 can be consumed from “/dev/input/mice”, “/dev/input/mouseX”, or “/dev/input/eventX”¹²⁸.

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 2.5.72¹²⁹.

```
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:     Vojtěch Pavlik <vojtech@suse.cz>
srcversion: 90003C5C18D7DE1706120
alias:      serio:tu0:pr*id=ex*
alias:      serio:tu0:pr*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:B0:9C:B0:D5:C0:8D:F7:05:30: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:0B:F3:E1:53:a5:4b:
            8F:01:F0:6B:F7:E7:a1:02:F4:6A:48:D0:43:FB:7B:3a:7E:0A:25:7B:
            35:85:60:CC:22:20:B9:4a:93:2a:FC:B0:02:36:0B:0a:0D:5E:C8:7B:
            BD:96:A6:66:72:D4:4C:87:64:59:E5:1D:76:0B:04:F2:20:70:93:D0:
            23:FF:BB:9F:57:F5:D2:0B:9C:98:B0:37:0B:05:39:72:38:50:87:7C:
            13:CS
parm:    tipdmg: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:    hpk_mode:default hpk 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:    reset_after_idle:device after sending hot packets (0 = never). (uint)
parm:    resume_time:How long can mouse stay idle before forcing resume (in seconds, 0 = never). (uint)
```

¹²⁷ <https://elixir.bootlin.com/linux/v6.2.6/source/drivers/input/mouse/psmouse-base.c#L2046>

¹²⁸ <https://www.kernel.org/doc/html/v5.5/input/input.html>

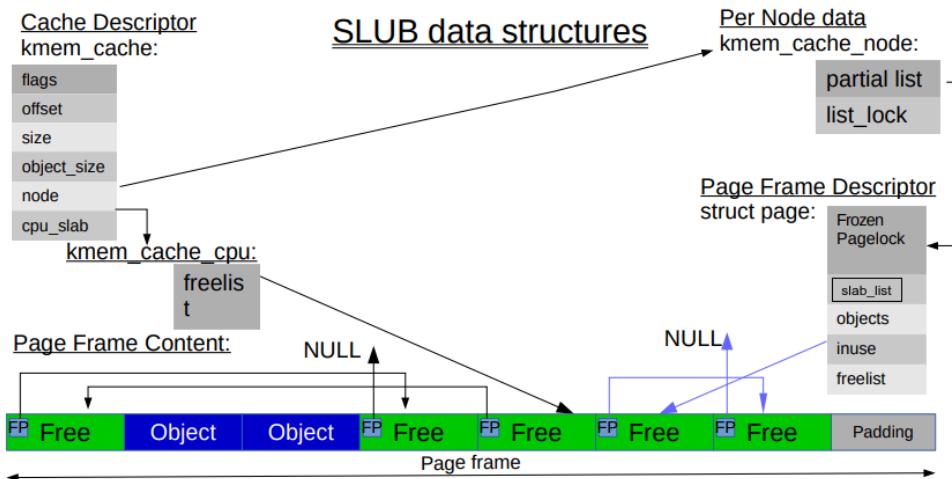
¹²⁹ <https://elixir.bootlin.com/linux/v2.5.72/source/drivers/input/mouse/psmouse-base.c>

Slub_flushwq (SLUB Flush Work Queue)

“slub_flushwq” is a kernel thread which based on a workqueue¹³⁰ 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 enabled¹³¹. 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 RAM¹³². The queuing of work is done inside the “flush_all_cpus_locked” function¹³³.

SLUB is also known as the “Unqueued Slab Allocator”¹³⁴. Slab allocation¹³⁵ 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/deallocations¹³⁶.

Thus, SLUB is a slab allocator that limits the use of cache lines instead of using queued object per cpu/per node list¹³⁷. 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 slab pages¹³⁸. 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 below¹³⁹.



¹³⁰ <https://elixir.bootlin.com/linux/v6.2.6/source/mm/slub.c#L5057>

¹³¹ <https://elixir.bootlin.com/linux/v6.2.6/source/mm/slub.c#L5056>

¹³² https://cateee.net/lkddb/web-lkddb/SLUB_TINY.html

¹³³ <https://elixir.bootlin.com/linux/v6.2.6/source/mm/slub.c#L2822>

¹³⁴ <https://lwn.net/Articles/229096/>

¹³⁵ <https://hammertux.github.io/slub-allocator>

¹³⁶ https://en.wikipedia.org/wiki/Slab_allocation

¹³⁷ <https://elixir.bootlin.com/linux/v6.2.6/source/mm/slub.c#L3>

¹³⁸ <https://hammertux.github.io/slub-allocator>

¹³⁹ <https://hammertux.github.io/img/SLUB-DS.png>

pgdatinit

“pgdatinit” is a kernel which is started by executing the “kthread_run()” function¹⁴⁰. The kernel thread executes the “deferred_init_memmap()” function¹⁴¹.

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]”¹⁴².

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 threads¹⁴³.

Lastly, after the initialization flow is finished an information message is sent to the kernel ring buffer¹⁴⁴ - as you can see in the image below¹⁴⁵.

```
[ 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
```

¹⁴⁰ https://elixir.bootlin.com/linux/v6.3-rc4/source/mm/page_alloc.c#L2284

¹⁴¹ https://elixir.bootlin.com/linux/v6.3-rc4/source/mm/page_alloc.c#L2108

¹⁴² https://elixir.bootlin.com/linux/v6.3-rc4/source/mm/page_alloc.c#L2283

¹⁴³ https://cateee.net/lkddb/web-lkddb/DEFERRED_STRUCT_PAGE_INIT.html

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

¹⁴⁵ <https://www.mail-archive.com/debian-bugs-dist@lists.debian.org/msg1822096.html>

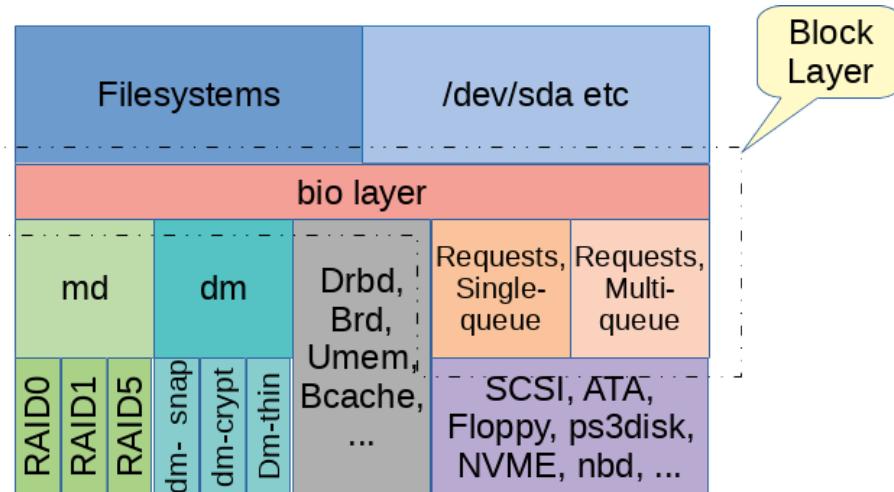
kblockd (Kernel Block Daemon)

“kblockd” is a kernel thread based on a workqueue¹⁴⁶ 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 below¹⁴⁷.

Overall, one might think that we can use keventd¹⁴⁸ 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 drivers¹⁴⁹.

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 devices¹⁵⁰.



¹⁴⁶ <https://elixir.bootlin.com/linux/v6.2.9/source/block/blk-core.c#L1191>

¹⁴⁷ <https://lwn.net/Articles/736534/>

¹⁴⁸ <https://lwn.net/Articles/11351/>

¹⁴⁹ <https://mirrors.edge.kernel.org/pub/linux/kernel/people/akpm/patches/2.5/2.5.70/2.5.70-mm8/broken-out/kblockd.patch>

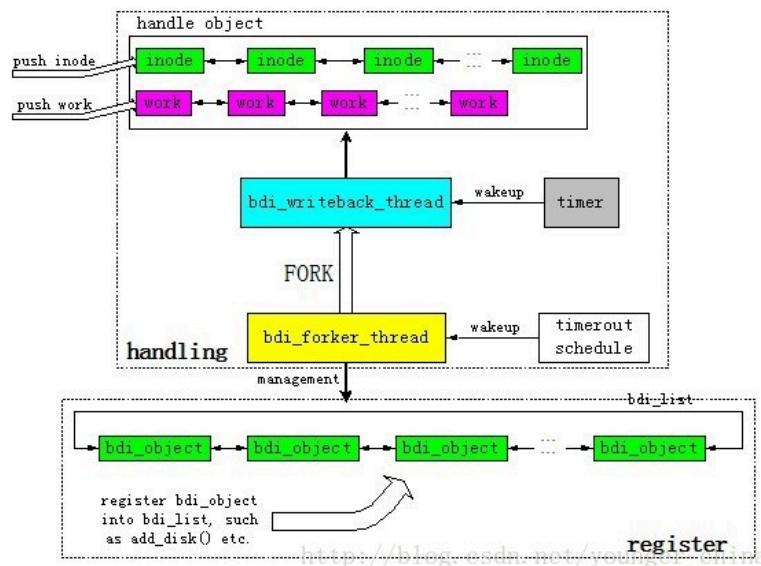
¹⁵⁰ <https://elixir.bootlin.com/linux/v6.3-rc4/source/block/blk-core.c#L13>

writeback

The kernel thread “writeback” is based on a workqueue¹⁵¹. The goal of the kernel thread is to serve all async writeback tasks¹⁵². 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 disk¹⁵³.

Overall, there 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 data¹⁵⁴.

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 disk¹⁵⁵. By the way, “bdi” stands for “Backing Device Information”¹⁵⁶. Lastly, to get an overview of the “writeback” mechanism you can checkout the diagram below¹⁵⁷.



¹⁵¹ <https://elixir.bootlin.com/linux/v6.2.5/source/mm/backing-dev.c#L363>

¹⁵² <https://elixir.bootlin.com/linux/v6.2.5/source/mm/backing-dev.c#L35>

¹⁵³ <https://www.oreilly.com/library/view/understanding-the-linux/0596005652/ch15s01.html>

¹⁵⁴ <https://lwn.net/Articles/326552/>

¹⁵⁵ https://blog.csdn.net/younger_china/article/details/55187057

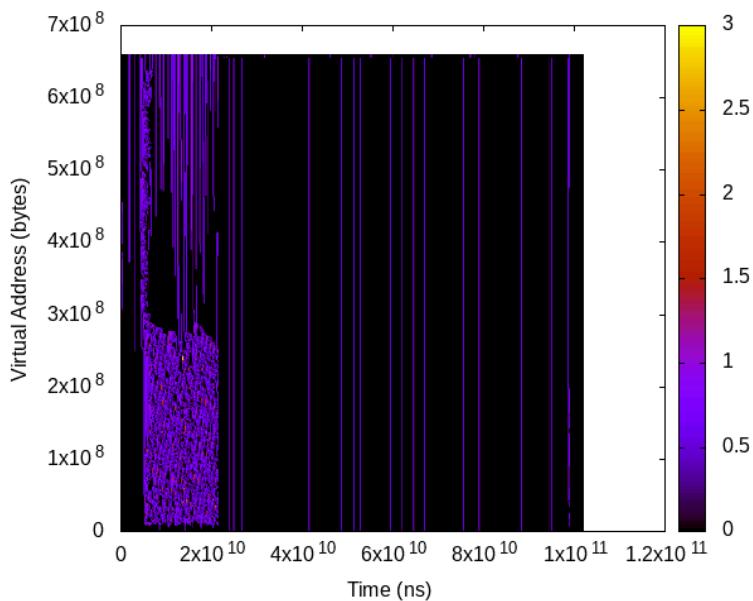
156 <https://lwn.net/Articles/326552/>157 https://blog.csdn.net/younger_china/article/details/55187057

kdamond (Data Access MONitor)

“kdamond” is a kernel thread which is created using the “kthread_run()” function¹⁵⁸ which is part of the DAMON (Data Access MONitor) subsystem. The kernel thread executes the “kdamon_fn()” function¹⁵⁹. Overall, DAMON provides a lightweight data access monitoring facility that can help users in analyzing the memory access patterns of their systems¹⁶⁰. Based on the documentation DAMON increases the memory usage by 0.12% and slows the workloads down by 1.39%¹⁶¹.

Also, DAMON has an API for kernel programs¹⁶². 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 configurations¹⁶³.

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 benchmarks¹⁶⁴. 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-blackscholes¹⁶⁵. Lastly, there are also other mechanisms in Linux that can help with data access monitoring such as “Perf Mem” and “Idle Page Tracking”



¹⁵⁸ <https://elixir.bootlin.com/linux/v6.3-rc5/source/mm/damon/core.c#L632>

¹⁵⁹ <https://elixir.bootlin.com/linux/v6.3-rc5/source/mm/damon/core.c#L1304>

¹⁶⁰ <https://www.kernel.org/doc/html/latest/admin-guide/mm/damon/index.html>

¹⁶¹ <https://damonitor.github.io/doc/html/v20/vm/damon/eval.html>

¹⁶² <https://www.kernel.org/doc/html/v5.17/vm/damon/api.html#functions>

¹⁶³ <https://sjn38.github.io/post/damon/>

¹⁶⁴ <https://parsec.cs.princeton.edu/parsec3-doc.htm>

¹⁶⁵ <https://lwn.net/Articles/813108/>

kintegrityd (Kernel Integrity Daemon)

“kintegrityd” is a kernel thread based on a workqueue¹⁶⁶ 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”¹⁶⁷. The function is called to complete a read request by verifying the transferred integrity metadata and then calls the original bio end_io function¹⁶⁸.

This procedure is done to ensure that the data was not changed by mistake (like in a case of a bug or an hardware failure¹⁶⁹). 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 layer¹⁷⁰.

Moreover, in order for it to work we should enable CONFIG_BLK_DEV_INTEGRITY, which is defined as “Block layer data integrity support”¹⁷¹. 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() function¹⁷². 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 userspace¹⁷³.

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

¹⁶⁶ <https://elixir.bootlin.com/linux/v6.1/source/block/bio-integrity.c#L455>

¹⁶⁷ <https://elixir.bootlin.com/linux/v6.1/source/block/bio-integrity.c#L317>

¹⁶⁸ <https://elixir.bootlin.com/linux/v6.1/source/block/bio-integrity.c#L313>

¹⁶⁹ <https://www.quora.com/What-is-the-purpose-of-kintegrityd-Linux-Kernel-Daemon/answer/Liran-Ben-Haim>

¹⁷⁰ <https://www.kernel.org/doc/Documentation/block/data-integrity.txt>

¹⁷¹ <https://elixir.bootlin.com/linux/v6.1/source/block/Kconfig#L60>

¹⁷² <https://www.kernel.org/doc/Documentation/block/data-integrity.txt>

¹⁷³ <https://www.kernel.org/doc/Documentation/block/data-integrity.txt>

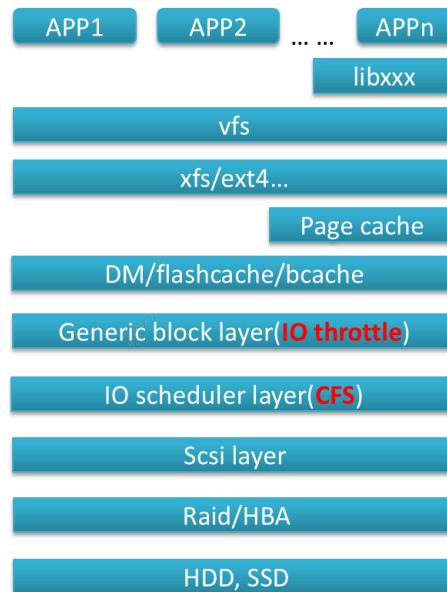
kthrotld (Kernel Throttling Daemon)

“kthrotld” is a kernel thread which was created using an workqueue¹⁷⁴ 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 queues¹⁷⁵.

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”¹⁷⁶ and then where is it referenced¹⁷⁷. 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 device¹⁷⁸.

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)¹⁷⁹.

Overall, IO throttling is done as part of the generic block layer and before the IO scheduler as seen in the diagram below¹⁸⁰. As described above this is also known as “Block Throttling”¹⁸¹.



¹⁷⁴ <https://elixir.bootlin.com/linux/v6.1.1/source/block/blk-throttle.c#L2470>

¹⁷⁵ <https://www.halolinux.us/kernel-architecture/request-queues.html>

¹⁷⁶ <https://elixir.bootlin.com/linux/v6.1.1/source/include/linux/blkdev.h#L395>

¹⁷⁷ https://elixir.bootlin.com/linux/v6.1.1/C/ident/request_queue

¹⁷⁸ <https://www.oreilly.com/library/view/linux-device-drivers/0596000081/ch12s04.html>

¹⁷⁹ <https://developer.aliyun.com/article/789736>

¹⁸⁰ <https://blog.csdn.net/yiveguzhou100/article/details/104044419>

¹⁸¹ <https://developer.aliyun.com/article/789736>

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>”¹⁸². It is the “SCSI error handler” which is responsible for all of the error handling targeting every SCSI host¹⁸³. The kernel thread is executing the “scsi_error_handler” function¹⁸⁴.

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 below¹⁸⁵.

Lastly, SCSI stands for “Small Computer System Interface”¹⁸⁶. It is a set of standards (from ANSI) for electronic interfaces in order to communicate with peripheral hardware like CD-ROM drives, tape drivers, printers, disk drives and more¹⁸⁷. For more information about SCSI I suggest going over.



© 2008 HowStuffWorks

¹⁸² <https://elixir.bootlin.com/linux/v6.4-rc1/source/drivers/scsi/hosts.c#L504>

¹⁸³ https://elixir.bootlin.com/linux/v6.4-rc1/source/drivers/scsi/scsi_error.c#L2230

¹⁸⁴ https://elixir.bootlin.com/linux/v6.4-rc1/source/drivers/scsi/scsi_error.c#L2233

¹⁸⁵ <https://computer.howstuffworks.com/scsi.htm>

¹⁸⁶ <https://hackaday.com/2023/03/02/scsi-the-disk-bus-for-everything/>

¹⁸⁷ <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” function¹⁸⁸. It is part of the common block controller cgroup interface¹⁸⁹.

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 cgroup¹⁹⁰. An example of priority inversion is shown in the diagram below¹⁹¹.

Thus, to avoid the problem mentioned above the function “submit_bio”¹⁹² punts the issuing of the bio request to a dedicated work item (per-block cgroup).

It calls “blkcg_punt_bio_submit”¹⁹³, which will call “__blkcg_punt_bio_submit”¹⁹⁴.

Priority inversion.

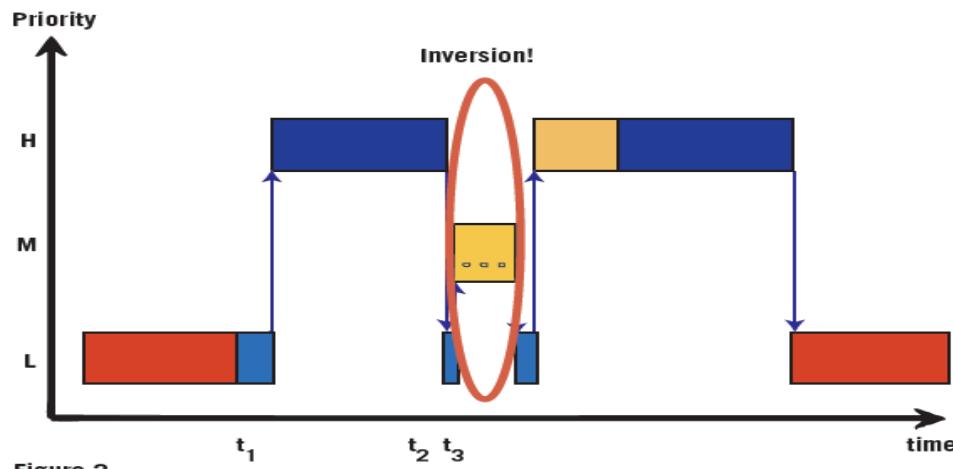


Figure 2

¹⁸⁸ <https://elixir.bootlin.com/linux/v6.2.5/source/block/blk-cgroup.c#L2058>

¹⁸⁹ <https://elixir.bootlin.com/linux/v6.2.5/source/block/blk-cgroup.c#L3>

¹⁹⁰ <https://patchwork.kernel.org/project/linux-block/patch/20190627203952.386785-6-tj@kernel.org/>

¹⁹¹ <https://embeddedgurus.com/barr-code/2010/11/firmware-specific-bug-8-priority-inversion/>

¹⁹² <https://elixir.bootlin.com/linux/v6.2.5/source/block/blk-core.c#L829>

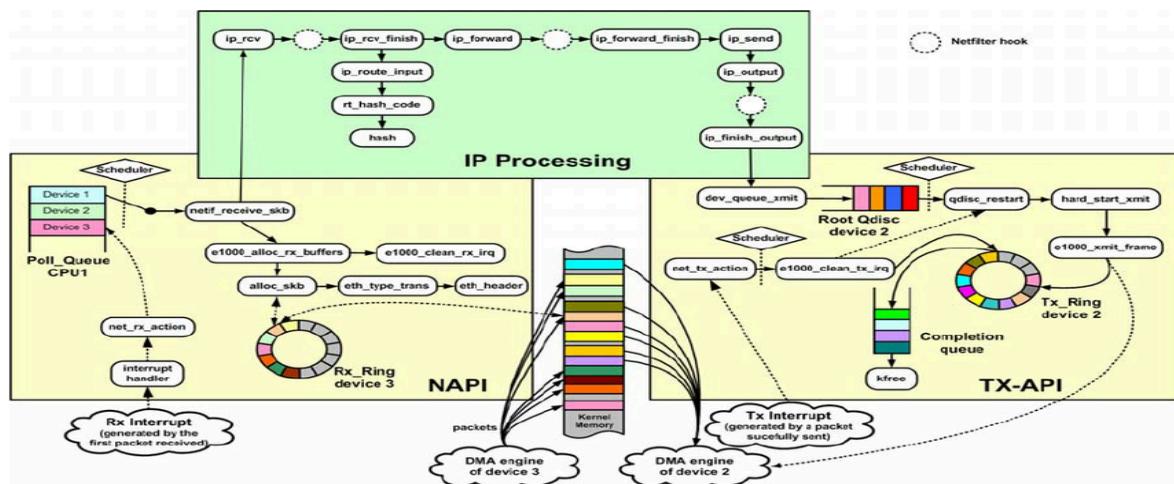
¹⁹³ <https://elixir.bootlin.com/linux/v6.2.5/source/block/blk-cgroup.h#L380>

¹⁹⁴ <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 pace¹⁹⁵. 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 compatibility¹⁹⁶. “napi” is a kernel thread which is created using the “kthread_run()”¹⁹⁷ 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”¹⁹⁸ 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 something¹⁹⁹. 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 again²⁰⁰. A detailed diagram of the NAPI flow is shown in the diagram below²⁰¹.



¹⁹⁵ <https://www.hitchhikersguidetolearning.com/2023/04/09/handling-receive-packets-via-napi/>

¹⁹⁶ <https://wiki.linuxfoundation.org/networking/napi>

¹⁹⁷ <https://elixir.bootlin.com/linux/v6.4-rc4/source/net/core/dev.c#L1371>

¹⁹⁸ <https://elixir.bootlin.com/linux/v6.4-rc4/source/net/core/dev.c#L662>

¹⁹⁹ <https://lwn.net/Articles/833840/>

²⁰⁰ <https://www.jianshu.com/p/7d4e36c0abe8>

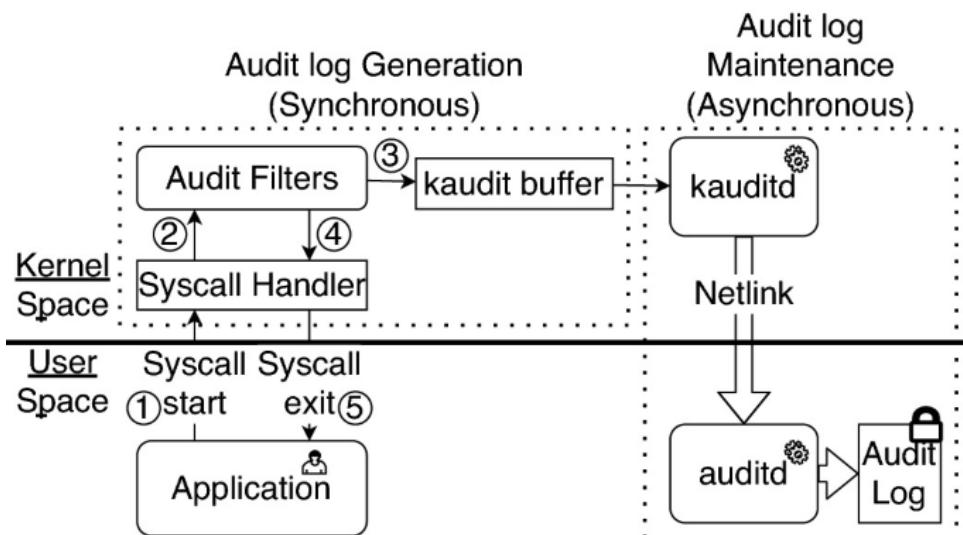
²⁰¹ <https://tinyurl.com/mrfnus8w>

kaudit (Kernel Audit Daemon)

“kaudit” is a kernel thread which is started using the “kthread_run” function²⁰². The kernel thread is calling the “kaudit_thread” function, this function is responsible for sending audit logs to userspace²⁰³. Overall, the kernel mechanism in the Linux kernel has a couple of goals: integrate fully with LSMs²⁰⁴, 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-mode²⁰⁵.

Thus, we can say “kaudit” is the kernel component of the “Linux Auditing System” which handles the audit events - as shown in the diagram below²⁰⁶. 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 rules²⁰⁷.

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 system²⁰⁸. Also, there are great examples for “audit.rules” in GitHub²⁰⁹.



²⁰² <https://elixir.bootlin.com/linux/v6.4-rc4/source/kernel/audit.c#L1700>

²⁰³ <https://elixir.bootlin.com/linux/v6.4-rc4/source/kernel/audit.c#L1828>

²⁰⁴ <https://medium.com/@boutnaru/linux-security-lsm-linux-security-modules-907bbcf8c8b4>

²⁰⁵ <https://elixir.bootlin.com/linux/v6.4-rc4/source/kernel/audit.c#L11>

²⁰⁶ https://link.springer.com/chapter/10.1007/978-3-031-17143-7_30

²⁰⁷ <https://manpages.debian.org/unstable/auditd/audit.rules.7.en.html>

²⁰⁸ <https://linux.die.net/man/8/auditctl>

²⁰⁹ <https://github.com/Neo23x0/auditd/blob/master/audit.rules>

tpm_dev_wq (Trusted Platform Module Device Work Queue)

“tpm_dev_wq” is a kernel thread base on a workqueue²¹⁰. It belongs a device file system interface for “Trusted Platform Module” aka TPM²¹¹.

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 storage²¹².

Moreover, a work is queued for “tpm_dev_wq” as part of the function “tpm_common_write”²¹³. In case we are working in non-blocking mode an async job for sending the command is scheduled²¹⁴.

Lastly, “tpm-dev-common.c” is compiled as part of the kernel TPM device drivers as shown in the Makefile²¹⁵. 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)
```

²¹⁰ <https://elixir.bootlin.com/linux/v6.3-rc7/source/drivers/char/tpm/tpm-dev-common.c#L273>

²¹¹ <https://elixir.bootlin.com/linux/v6.3-rc7/source/drivers/char/tpm/tpm-dev-common.c#L13>

²¹² https://wiki.archlinux.org/title/Trusted_Platform_Module

²¹³ <https://elixir.bootlin.com/linux/v6.3-rc7/source/drivers/char/tpm/tpm-dev-common.c#L209>

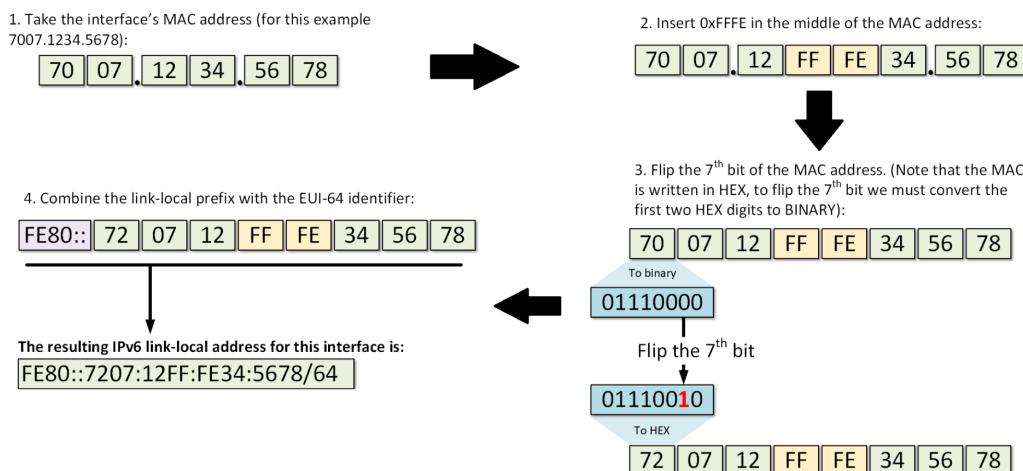
²¹⁴ <https://elixir.bootlin.com/linux/v6.3-rc7/source/drivers/char/tpm/tpm-dev-common.c#L202>

²¹⁵ <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 workqueue²¹⁶. This code is part of the Linux INET6 implementation and is responsible for the IPv6 Address auto configuration²¹⁷. 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 assignment²¹⁸.

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`”²¹⁹. We can see there the link-local prefix²²⁰ and the call for generating the EUI-64 identifier by the function “`ipv6_generate_eui64`”²²¹. 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 more²²². 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`”²²³.

²¹⁶ <https://elixir.bootlin.com/linux/v6.2.11/source/net/ipv6/addrconf.c#L7292>

²¹⁷ <https://elixir.bootlin.com/linux/v6.2.11/source/net/ipv6/addrconf.c#L13>

²¹⁸ <https://www.networkacademy.io/crna/ipv6/stateless-address-autoconfiguration-slaac>

²¹⁹ <https://elixir.bootlin.com/linux/v6.2.11/source/net/ipv6/addrconf.c#L3314>

²²⁰ <https://elixir.bootlin.com/linux/v6.2.11/source/net/ipv6/addrconf.c#L3326>

²²¹ <https://elixir.bootlin.com/linux/v6.2.11/source/net/ipv6/addrconf.c#L3345>

²²² <https://datatracker.ietf.org/doc/html/rfc4862>

²²³ <https://elixir.bootlin.com/linux/v6.2.11/source/net/ipv6/addrconf.c#L4058>

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.

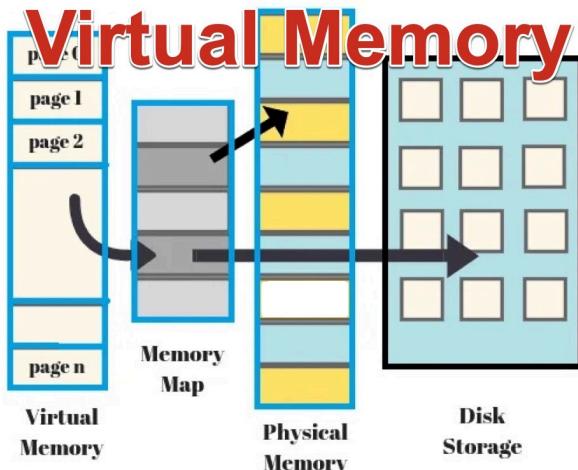
mm_percpu_wq (Per-CPU Memory Work Queue)

“mm_percpu_wq” is a kernel thread based on a workqueue which is created in the “init_mm_internals” function²²⁴. It is part of the the statistics management regarding virtual memory²²⁵. An overview diagram of virtual memory is shown below²²⁶.

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”²²⁷. “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”²²⁸. This proc entry is created with others (“buddyinfo”, “pagetypeinfo” and “zoneinfo”) in the same file in which “mm_percpu_mm” is allocated²²⁹. We can see the list of the metric counters in the source code²³⁰.

As it names suggested the kernel thread is responsible for accumulating the vm events among all CPUs²³¹. It is done by going over all the “online” CPUs²³². Lastly, we can use different cli tools to review the different statistic counters. One of those tools is “vmstat”²³³.



²²⁴ <https://elixir.bootlin.com/linux/v6.4-rc5/source/mm/vmstat.c#L2100>

²²⁵ <https://elixir.bootlin.com/linux/v6.4-rc5/source/mm/vmstat.c#L5>

²²⁶ <https://iboysoft.com/wiki/virtual-memory.html>

²²⁷ <https://elixir.bootlin.com/linux/v6.4-rc5/source/mm/vmstat.c#L2021>

²²⁸ <https://man7.org/linux/man-pages/man5/proc.5.html>

²²⁹ <https://elixir.bootlin.com/linux/v6.4-rc5/source/mm/vmstat.c#L2123>

²³⁰ <https://elixir.bootlin.com/linux/v6.4-rc5/source/mm/vmstat.c#L1168>

²³¹ <https://elixir.bootlin.com/linux/v6.4-rc5/source/mm/vmstat.c#L126>

²³² <https://elixir.bootlin.com/linux/v6.4-rc5/source/mm/vmstat.c#L117>

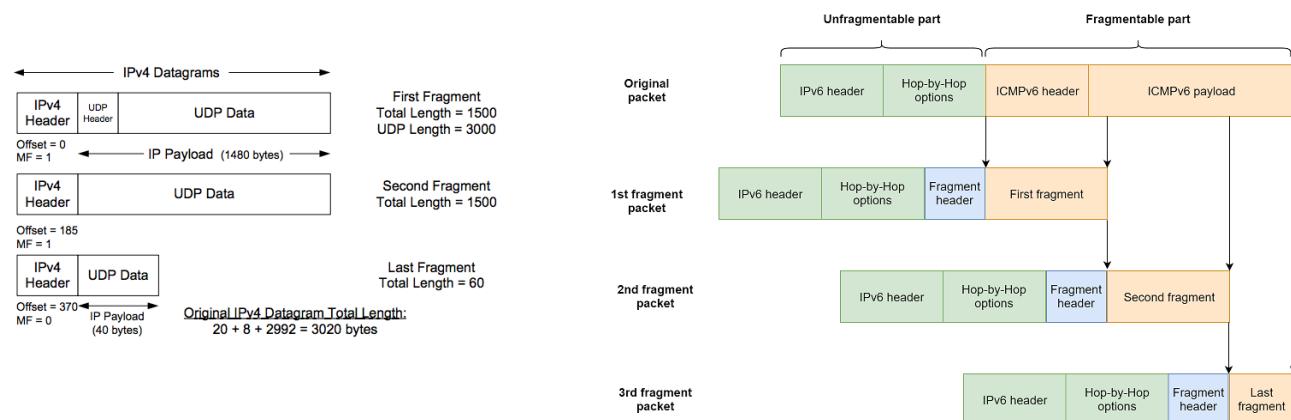
²³³ <https://linux.die.net/man/8/vmstat>

inet_frag_wq (IP Fragmentation Work Queue)

The kernel thread “inet_frag_wq” is created using a workqueue²³⁴, we could have guessed it based on a workqueue due 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 packets²³⁵.

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”²³⁶. An illustration of the flow is shown in the diagram below both for IPv4²³⁷ and IPv6²³⁸.

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"²³⁹. 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 fragments²⁴⁰.



²³⁴ https://elixir.bootlin.com/linux/v6.2-rc1/source/net/ipv4/inet_fragment.c#L211

²³⁵ https://elixir.bootlin.com/linux/v6.2-rc1/source/net/ipv4/inet_fragment.c#L6

²³⁶ <https://www.geeksforgeeks.org/ipv6-fragmentation-header/>

²³⁷ <https://notes.shichao.io/tcpv1/ch10/>

²³⁸ <https://blog.quarkslab.com/analysis-of-a-windows-ipv6-fragmentation-vulnerability-cve-2021-24086.html>

²³⁹ https://elixir.bootlin.com/linux/v6.2-rc1/source/include/net/inet_frag.h#L66

²⁴⁰ https://elixir.bootlin.com/linux/v6.2-rc1/source/net/ipv4/in_fragment.c#L411

kstrp (Stream Parser)

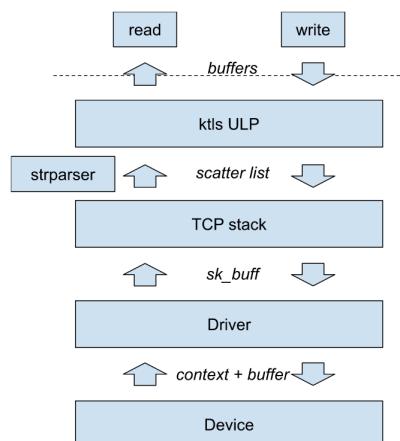
“kstrp” is based on a single threaded workqueue²⁴¹. From the source code documentation we can see that “strparser” means “Stream Parser”²⁴². A stream parser is a utility that gets data streams and parses the application layer protocol over those streams. A stream parser can work in one of two modes: 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 socket²⁴³.

Thus, we can say that we can parse application layer protocol messages in TCP. It is basically a generalization of KCM (Kernel Connection Multiplexor)²⁴⁴.

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 TCP²⁴⁵.

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 module²⁴⁶. One example for that is the implementation of KTLS²⁴⁷ (a Linux TLS/DTLS kernel module). An illustration of the flow is shown below²⁴⁸.



²⁴¹ <https://elixir.bootlin.com/linux/v6.1.1/source/net/strparser/strparser.c#L539>

²⁴² <https://elixir.bootlin.com/linux/v6.1.1/source/net/strparser/strparser.c#L3>

²⁴³ <https://www.kernel.org/doc/html/v5.10/networking/strparser.html>

²⁴⁴ <https://lwn.net/Articles/695982/>

²⁴⁵ <https://www.kernel.org/doc/html/latest/networking/kcm.html>

²⁴⁶ <https://zhuanlan.zhihu.com/p/543663512>

²⁴⁷ https://github.com/ktls/af_ktls

²⁴⁸ <https://docs.kernel.org/networking/tls-offload.html>

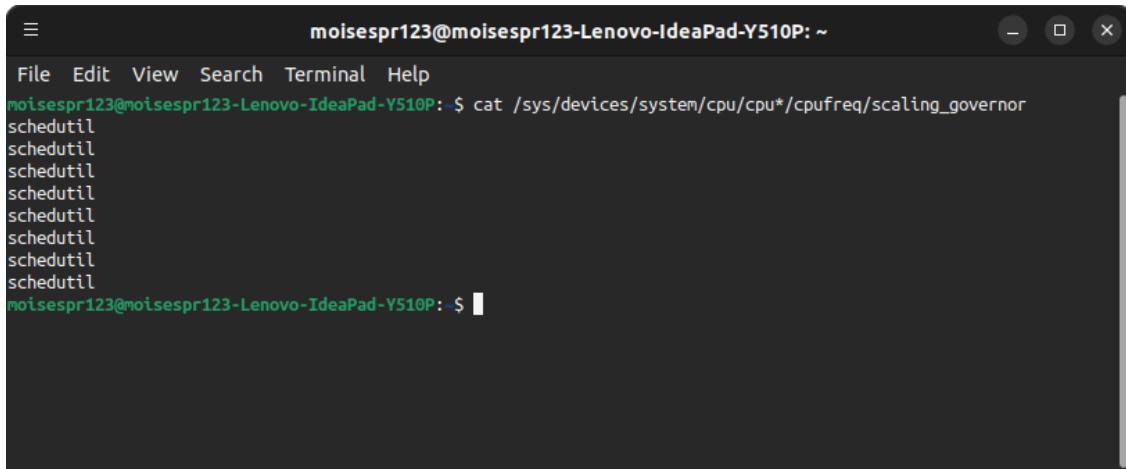
devfreq_wq

“devfreq_wq” is a kernel thread which is based on a freezable workqueue²⁴⁹. It is part of the Generic Dynamic Voltage and Frequency Scaling (DVFS) Framework for Non-CPU devices²⁵⁰.

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 time²⁵¹.

Moreover, based on the Linux documentation there are 6 governors: “Performance”, “Powersave”, “Userspace”, “Ondemand”, “Conservative” and “Schedutil”²⁵². We can also develop our own governor as a kernel module, we just need to register it using the function “cpufreq_register_governor”²⁵³.

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 more²⁵⁴. An example of reading the current configure governor is shown below²⁵⁵.



```
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
moisespr123@moisespr123-Lenovo-IdeaPad-Y510P:~$
```

²⁴⁹ <https://elixir.bootlin.com/linux/v6.2.5/source/drivers/devfreq/devfreq.c#L1997>

²⁵⁰ <https://elixir.bootlin.com/linux/v6.2.5/source/drivers/devfreq/devfreq.c#L3>

²⁵¹ [https://wiki.somlabs.com/index.php/How to scale CPU frequency with DVFS framework](https://wiki.somlabs.com/index.php/How_to_scale_CPU_frequency_with_DVFS_framework)

²⁵² <https://www.kernel.org/doc/Documentation/cpu-freq/governors.txt>

²⁵³ <https://elixir.bootlin.com/linux/v6.5-rc2/source/drivers/cpufreq/cpufreq.c#L2443>

254 <https://www.kernel.org/doc/Documentation/cpu-freq/user-guide.txt>255 <https://moisescardona.me/changing-the-cpu-governor-to-performance-in-linux/>

dmcrypt_write (Device Mapper for Transparent Encryption/Decryption)

“dmcrypt_write” is a kernel thread which is created using the “kthread_run” function²⁵⁶. The name of the kernel thread is in the pattern of “dmcrypt_write/%s”, where the added string represents the device name. Overall, “dm-crypt” is a device-mapper target²⁵⁷ supported from kernel version 2.6.4²⁵⁸. It is responsible for transparent (aka real-time/on-the-fly encryption) block device encryption while using the kernel crypto API²⁵⁹.

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 kernel²⁶⁰. Moreover, the function that is executed as part of the kernel thread is “dmcrypt_write” function²⁶¹. 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 dmcrypt_write
0000000000000000 t dmcrypt_write           [dm_crypt]
```

²⁵⁶ <https://elixir.bootlin.com/linux/v6.5-rc3/source/drivers/md/dm-crypt.c#L3388>

²⁵⁷ <https://elixir.bootlin.com/linux/v6.5-rc3/source/drivers/md/dm-crypt.c#L3689>

²⁵⁸ <https://elixir.bootlin.com/linux/v2.6.4/source/drivers/md/dm-crypt.c>

²⁵⁹ <https://gitlab.com/cryptsetup/cryptsetup/-/wikis/DMCrypt>

²⁶⁰ <https://elixir.bootlin.com/linux/v6.5-rc3/source/drivers/md/Makefile#L59>

²⁶¹ <https://elixir.bootlin.com/linux/v6.5-rc3/source/drivers/md/dm-crypt.c#L1922>

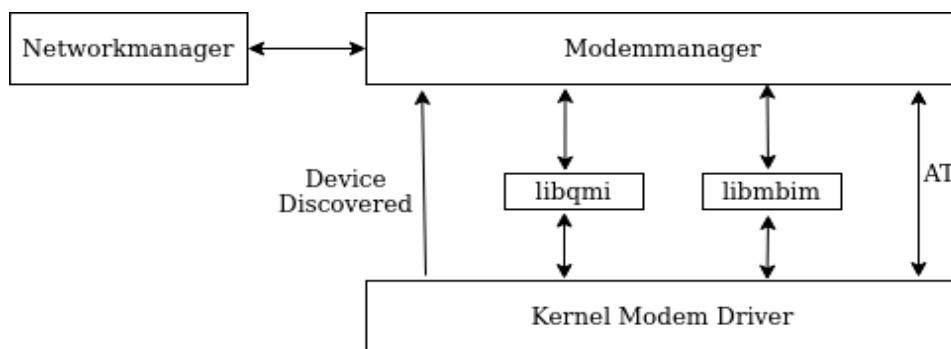
ModemManager (Modem Management Daemon)

“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 modems²⁶². Also, it is started by PID 1 (init/systemd) with the permission of the root user.

Overall, it is a DBus-powered²⁶³ 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 more²⁶⁴.

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 OpenWRT²⁶⁵.

It is important to understand that “ModemManager” leverages “libqmi”²⁶⁶ and “libmbim”²⁶⁷ to communicate over QMI (Qualcomm MSM Interface) and MBIM (Mobile Interface Broadband Model) for setting connection to the cellular network²⁶⁸. 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 repo²⁶⁹. I also suggest going over the documentation site of “ModemManager” and the relevant libraries: libmbim, libqmi and libqrtr-glib²⁷⁰.



²⁶² <https://manpages.ubuntu.com/manpages/trusty/man8/ModemManager.8.html>

²⁶³ <https://www.freedesktop.org/software/ModemManager/api/latest/>

²⁶⁴ <https://manpages.ubuntu.com/manpages/trusty/man8/mmcli.8.html>

²⁶⁵ <https://modemmanager.org/>

²⁶⁶ <https://github.com/linux-mobile-broadband/libqmi>

²⁶⁷ <https://github.com/linux-mobile-broadband/libmbim>

²⁶⁸ <https://developer.toradex.com/software/connectivity/modem-support/>

²⁶⁹ <https://gitlab.freedesktop.org/mobile-broadband/ModemManager>

²⁷⁰ <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.org²⁷¹. An example of such oops is shown in the screenshot below²⁷². By the way, they are also known as “soft panic”²⁷³.

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 problem²⁷⁴.

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 information²⁷⁵. Also, I recommend also enabling “CONFIG_FRAME_POINTER”, which gives very useful debugging information in case of kernel bugs - precise oopses/stacktraces/warnings²⁷⁶.

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 10000²⁷⁷.

```
(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:fffffb41340e6fdd8 EFLAGS: 00010246
(12710.156849) RAX: 0000000000000019 RBX: ffffffff81015040 RCX: 0000000000000000
(12710.157513) RDX: 0000000000000000 RSI: ffffffff83bc9d39 RDI: 00000000fffffff
(12710.158171) RBP: fffff8d6101bd1d50 R08: 0000000000000000 R09: fffffb41340e6fc90
(12710.158826) R10: 0000000000000003 R11: ffffff83f3d1e8 R12: fffffb41340e6fde0
(12710.159483) R13: 0000000000000000 R14: 0000000000000000 R15: 0000000000000000
(12710.160143) FS: 00007f6c290b31c0(0000) GS:ffff8d6411a00000(0000) knlGS:00000
(12710.160820) CS: 0010 DS: 0000 ES: 0000 CR0: 0000000080050033
(12710.161478) CR2: 0000000004134f0 CR3: 0000000018be34005 CR4: 000000000003706f0
(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
```

²⁷¹ <https://linux.die.net/man/8/kerneloops>

²⁷² <https://nakedsecurity.sophos.com/2023/03/13/linux-gets-double-quick-double-update-to-fix-kernel-oops/>

²⁷³ <https://www.opensourceforu.com/2011/01/understanding-a-kernel-oops/>

²⁷⁴ https://en.wikipedia.org/wiki/Linux_kernel_oops

²⁷⁵ <https://www.oreilly.com/library/view/linux-device-drivers/0596005903/ch04.html>

²⁷⁶ https://cateee.net/lkddb/web-lkddb/FRAME_POINTER.html

²⁷⁷ <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 them²⁷⁸ - 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 xargs²⁷⁹.

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 usage²⁸⁰. Lastly, there are different implementations for xargs two examples are Apple’s one²⁸¹ and the busybox one²⁸².

```
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:~# _
```

²⁷⁸ <https://www.wikiwand.com/en/Xargs>

²⁷⁹ <https://man7.org/linux/man-pages/man1/xargs.1.html>

²⁸⁰ <https://www.tecmint.com/xargs-command-examples/>

²⁸¹ https://opensource.apple.com/source/shell_cmds/shell_cmds-149/xargs/

²⁸² <https://github.com/josefbsck/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 contracts. We use it when writing code in C/C++/Objective-C²⁸³.

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 phase²⁸⁴. This is do using “ifdef”, “#ifndef”, “#define”, “#else” and “#endif”²⁸⁵.

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 "Tr0LLer"

#ifndef 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 []="Tr0LLer";
    return;
}
```

²⁸³ <https://linux.die.net/man/1/cpp>

²⁸⁴ <https://www.programiz.com/c-programming/c-preprocessor-macros>

²⁸⁵ <https://www.cprogramming.com/reference/preprocessor/ifndef.html>

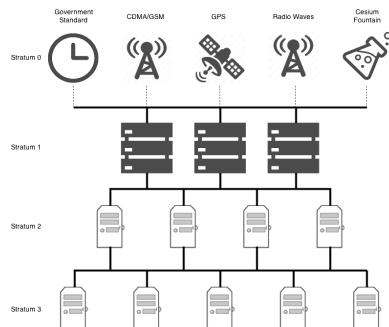
ntpd (Network Time Protocol Daemon)

“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 picoseconds²⁸⁶.

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 sources²⁸⁷. 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 123²⁸⁸.

Moreover, there is an open source implementation of NTP from the University of Delaware²⁸⁹. You can also check out a port of it for the Windows operating system²⁹⁰. There are two utilities that can be used when “ntpd” is running for tasks like troubleshooting/configuration/monitoring: “ntpq”²⁹¹ “ntpmon”²⁹². 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 below²⁹³.



²⁸⁶ <https://linux.die.net/man/8/ntpd>

²⁸⁷ <https://www.mankier.com/8/ntpd>

²⁸⁸ <https://www.techtarget.com/searchnetworking/definition/Network-Time-Protocol>

²⁸⁹ <http://www.ntp.org/>

²⁹⁰ <https://www.meinbergglobal.com/english/sw/ntp.htm>

²⁹¹ <https://www.mankier.com/1/ntp>

²⁹² <https://www.mankier.com/1/ntpmon>

²⁹³ [https://chrissshort.net/ntp-i-need-you-to-go-ahead-and-love-it/](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++ applications²⁹⁴.

Moreover, “gold” today is part of GNU binutils²⁹⁵. 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 layer²⁹⁶.

Thus, we can specify “gold” as the the linker to use by setting LD as part of a makefile, setting the LD environment variable²⁹⁷ or by using the “`-fuse-ld=gold`” option of gcc²⁹⁸.

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”²⁹⁹. “gold” (sometimes also “`ld.gold`”) can take different command line arguments which can affect the linking processes³⁰⁰.

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.

²⁹⁴ <https://opensource.googleblog.com/2008/04/gold-google-releases-new-and-improved.html>

²⁹⁵ <https://sourceware.org/binutils/>

²⁹⁶ [https://en.wikipedia.org/wiki/Gold_\(linker\)](https://en.wikipedia.org/wiki/Gold_(linker))

²⁹⁷ [https://en.wikipedia.org/wiki/Gold_\(linker\)](https://en.wikipedia.org/wiki/Gold_(linker))

²⁹⁸ <https://man7.org/linux/man-pages/man1/gcc.1.html>

²⁹⁹ <https://static.googleusercontent.com/media/research.google.com/en/pubs/archive/34417.pdf>

³⁰⁰ <https://www.mankier.com/1/ld.gold>

strace (System Call Tracing)

“strace” (System Call Tracing) is an ELF file which is located at “/bin/strace”, by the way “/bin” is usually a symbolic link to “/usr/bin” so we can find the binary at “/usr/bin/strace”. It can be leveraged for: tracing\monitoring\diagnosing\debugging Linux’s tasks (processes\threads)³⁰¹. It is used for displaying information about system calls³⁰² executed by a specific application. It is important to understand that strace is dependent on the “ptrace” syscall for performing its job³⁰³.

Overall, by default every line in the output of the strace command (which writes to standard error) contains: the system call name, its arguments in parentheses and its return value of the syscall. In case the syscall is not identified by strace it is displayed as “syscall_0x{NUM}”, where {NUM} is the syscall number in hex. Thus, we can say strace is used to trace system calls and signals until the monitored application exits³⁰⁴.

Lastly, it is important to know that using the different switches of strace we can use more advanced features such as (but not limited to): attach to an already running process (-p), print paths and more information associated with file descriptors (-y), filter by type of syscall (-e trace={SYSCALL}), trace only syscalls accessing a specific path (-P), perform syscall fault injection (-e fault={SYSCALL}), dump all data read\written to file descriptors in hex\ascii (-ewrite=1) and count time\calls\errors for every syscall (-c) - as shown in the screenshot below³⁰⁵. For more information I suggest going over the source code as part of strace’s GitHub repository³⁰⁶.

\$ strace -c ls > /dev/null	% time	seconds	usecs/call	calls	errors	syscall
	89.76	0.008016	4	1912		getdents
	8.71	0.000778	0	11778		lstat
	0.81	0.000072	0	8894		write
	0.60	0.000054	0	943		open
	0.11	0.000010	0	942		close
	0.00	0.000000	0	1		read
	0.00	0.000000	0	944		fstat
	0.00	0.000000	0	8		mmap
	0.00	0.000000	0	4		mprotect
	0.00	0.000000	0	1		munmap
	0.00	0.000000	0	7		brk
	0.00	0.000000	0	3	3	access
	0.00	0.000000	0	1		execve
	0.00	0.000000	0	1		sysinfo
	0.00	0.000000	0	1		arch_prctl
	100.00	0.008930		25440	3	total

³⁰¹ <https://www.geeksforgeeks.org/linux-unix-strace-command-in-linux-with-examples/>

³⁰² <https://medium.com/@boutinari/the-linux-concept-journey-syscalls-system-calls-efcd7703e072>

³⁰³ <https://man7.org/linux/man-pages/man2/ptrace.2.html>

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

³⁰⁵ <https://strace.io/>

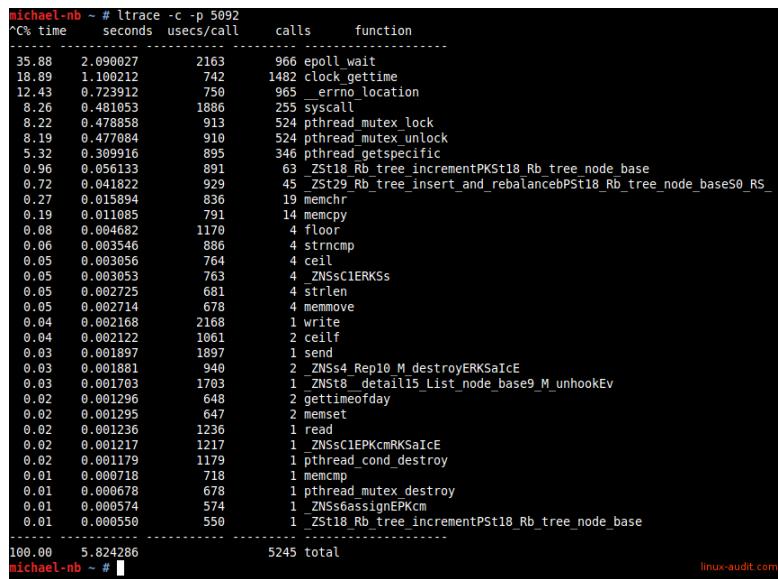
³⁰⁶ <https://github.com/strace/strace>

Itrace (Library Call Tracer)

“Itrace” (Library Call Tracing) is an ELF file which is located at “/bin/ltrace”, by the way “/bin” is usually a symbolic link to “/usr/bin” so we can find the binary at “/usr/bin/ltrace”. It is a CLI utility which can be used to monitor\trace a specific program until it exists. Its main goal is to tap and log the dynamic library calls which are invoked by the executed process. Also, it can record the signals which are received by that process. By the way, it can also display syscalls as well as library calls using the “-S” switch³⁰⁷. We can also filter the output of the command based on specific library functions and display statistics on the function used by a traced task - as shown in the screenshot below³⁰⁸.

Overall, as with strace utility³⁰⁹ also ltrace leverages the “ptrace” syscall³¹⁰ in order to set breakpoints in the task (process\thread) being traced³¹¹. From the source code of ltrace we deduct the following flow: starting a process\attaching to a running process, setting breakpoint(s) on library calls\syscalls and decoding library\syscall name and arguments³¹².

Lastly, it is important to understand that ltrace is dependent that the binary being traced is compiled and linked with lazy loading (gcc -z lazy {SOURCE_FILE} -o {OUTPUT_FILE}), if not the library information is not displayed³¹³.



```
michael-nb ~ # ltrace -c -p 5092
^% time    seconds  usecs/call   calls   function
-----
35.88  2.090027   2163      966 epoll_wait
18.89  1.100212    742     1482 clock_gettime
12.43  0.723912    750      965 __errno_location
 8.26  0.481053   1886      255 syscall
 8.22  0.478858    913      524 pthread_mutex_lock
 8.19  0.477084    910      524 pthread_mutex_unlock
 5.32  0.309916    895      346 pthread_getspecific
 0.96  0.056133    891      63 __ZSt18_Rb_tree_incrementPKSt18_Rb_tree_node_base
 0.72  0.041822    929      45 __ZSt29_Rb_tree_insert_and_rebalancebPSt18_Rb_tree_node_base50_RS_
 0.27  0.015894    836      19 memchr
 0.19  0.011085    791      14 memcpy
 0.08  0.004682   1170      4 floor
 0.06  0.003546    886      4 strcmp
 0.05  0.003056    764      4 ceil
 0.05  0.003053    763      4 __ZNssC1ERKSS
 0.05  0.002725    681      4 strlen
 0.05  0.002714    678      4 memmove
 0.04  0.002168   2168      1 write
 0.04  0.002122   1061      2 ceilf
 0.03  0.001897   1897      1 send
 0.03  0.001881    940      2 __ZNss4_Repl10_M_destroyERKSaIcE
 0.03  0.001703   1703      1 __ZNSt8_detail15_List_node_base9_M_unhookEv
 0.02  0.001296    648      2 gettimeofday
 0.02  0.001295    647      2 memset
 0.02  0.001236   1236      1 read
 0.02  0.001217   1217      1 __ZNssC1EPKcmRKSaIcE
 0.02  0.001179   1179      1 pthread_cond_destroy
 0.01  0.000718    718      1 memcmp
 0.01  0.000678    678      1 pthread_mutex_destroy
 0.01  0.000574    574      1 __ZNs6assignEPKcm
 0.01  0.000550    550      1 __ZSt18_Rb_tree_incrementPKSt18_Rb_tree_node_base
100.00  5.824286           5245 total
michael-nb ~ #
```

³⁰⁷ <https://man7.org/linux/man-pages/man1/ltrace.1.html>

³⁰⁸ <https://linux-audit.com/monitor-file-access-by-linux-processes/>

³⁰⁹ <https://medium.com/@boutnaru/the-linux-concept-journey-strace-system-call-tracing-11354f94b00d>

³¹⁰ <https://linux.die.net/man/2/ptrace>

³¹¹ <https://www.kernel.org/doc/ols/2007/ols2007v1-pages-41-52.pdf>

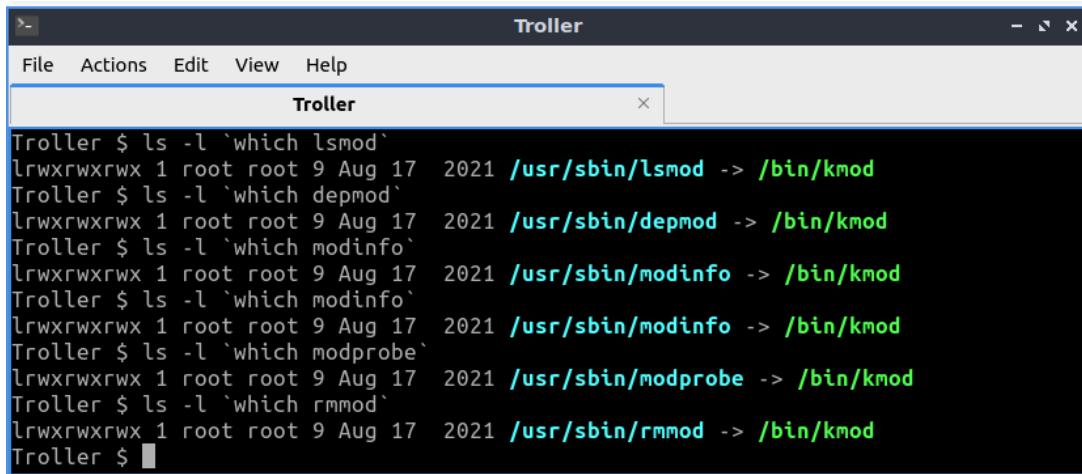
³¹² <https://gitlab.com/cespedes/ltrace>

³¹³ <https://medium.com/@boutnaru/linux-why-ltrace-does-not-work-on-new-versions-of-ubuntu-b2d89b55b70f>

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 modules³¹⁴. 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 directly³¹⁵.

Overall, by using “kmod” we can load/remove/insert/show information/resolve dependencies of kernel modules³¹⁶. “kmod” is developed by Lucas De Marchi³¹⁷. Lastly, the most well known utilities used to control kernel modules are just symbolic links to “kmod” examples are: lsmod, rmmod, 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 rmmod`
lrwxrwxrwx 1 root root 9 Aug 17 2021 /usr/sbin/rmmod -> /bin/kmod
Troller $
```

³¹⁴ https://linux-kernel-labs.github.io/refs/heads/master/labs/kernel_modules.html

³¹⁵ <https://man7.org/linux/man-pages/man8/kmod.8.html>

³¹⁶ <https://linuxhint.com/linux-kmod-command/>

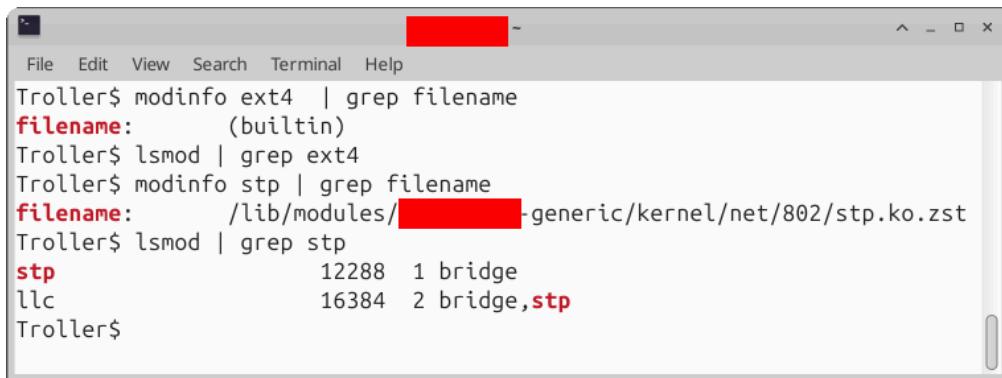
³¹⁷ <https://github.com/kmod-project/kmod>

lsmod (List Kernel Modules)

“lsmod” is a symbolic link³¹⁸ to “/sbin/lsmod” or “/usr/sbin/lsmod”. It references the ELF binary “kmod” which is located at “/bin/kmod” or “/usr/bin/kmod”³¹⁹. lsmod is used for showing the status of LKMs³²⁰. This is done by parsing the content of “/proc/modules”³²¹.

Overall, it is important to understand that in case of builtin kernel modules³²² which are compiled as part of the kernel itself, they won’t appear in the output of lsmod - as shown in the screenshot below.

Lastly, lsmod also reads information from “/sys/module/*”. For example in the function “kmod_module_get_size”³²³ which is called from “kmod_list_FOREACH”³²⁴. The output of lsmod provides the following data: “Module” (the name of the module), “Size” (in bytes and not memory in use) and “Used By” (the number of times it is used by running programs). Also, the list on the right is the name(s) of the module(s) which refer to the current one described in the specific row³²⁵ - as shown below. It is important to understand that lsmod could be incomplete in regards to the used by\reference counters³²⁶.



```
File Edit View Search Terminal Help
Troller$ modinfo ext4 | grep filename
filename:          (builtin)
Troller$ lsmod | grep ext4
Troller$ modinfo stp | grep filename
filename:          /lib/modules/[REDACTED]-generic/kernel/net/802/stp.ko.zst
Troller$ lsmod | grep stp
stp                  12288  1 bridge
llc                   16384  2 bridge,stp
Troller$
```

³¹⁸ <https://medium.com/@boutnaru/the-linux-concept-journey-linux-file-types-4cb622887331>

³¹⁹ <https://medium.com/@boutnaru/the-linux-process-journey-kmod-linux-kernel-module-handling-8ec844216436>

³²⁰ <https://medium.com/@boutnaru/the-linux-concept-journey-loadable-kernel-module-lkm-5eaad4db346a1>

³²¹ <https://man7.org/linux/man-pages/man8/lsmod.8.html>

³²² <https://medium.com/@boutnaru/the-linux-concept-journey-built-in-kernel-modules-0b9fd2bc524c>

³²³ <https://github.com/kmod-project/kmod/blob/master/libkmod/libkmod-module.c#L1482>

³²⁴ <https://github.com/kmod-project/kmod/blob/master/tools/lsmod.c#L98C2-L98C19>

³²⁵ <https://en.wikipedia.org/wiki/Lsmod>

³²⁶ <https://stackoverflow.com/questions/4073152/how-to-get-complete-dependency-list-of-kernel-modules-at-runtime/4238010#4238010>

insmod (Insert Kernel Module)

“insmod” is a symbolic link³²⁷ to that is located at “/sbin/insmod” or “/usr/sbin/insmod”. It references the ELF binary “kmod” which is located at “/bin/kmod” or “/usr/bin/kmod”³²⁸. insmod is used for inserting a LKM³²⁹ into the Linux kernel. We can use the “-f” (“--force”) switch of “insmod” for ignoring the module version and vermagic fields when loading and “-s” (“--syslog”) for sending errors to syslog instead of standard error³³⁰.

Overall, in order to load a kernel module we can use the “init_module” syscall or the “finit_module” syscall. The difference between the two is that “finit_module” reads the module to be loaded from a file descriptor (fd) while “init_module” gets a pointer to a buffer containing the binary image to be loaded³³¹. The current version of insmod uses the “finit_module” syscall in order to load a LKM - as shown in the screenshot below (the syscall fails because the kernel module is not signed). I suggest going over the code for more details³³².

Lastly, it is important to understand that insmod is intended to insert\load a single kernel module from any location. This is due to the fact it does not manage dependencies and does not read modules from a specific location³³³. In case of any errors it's best to check out “dmesg”³³⁴. By the way, insmod causes the module's init function of the LKM³³⁵.

% time	seconds	usecs/call	calls	errors	syscall
100.00	0.000152	152	1	1	finit_module
100.00	0.000152	152	1	1	total

³²⁷ <https://medium.com/@boutnaru/the-linux-concept-journey-linux-file-types-4cb622887331>

³²⁸ <https://medium.com/@boutnaru/the-linux-process-journey-kmod-linux-kernel-module-handling-8ec844216436>

³²⁹ <https://medium.com/@boutnaru/the-linux-concept-journey-loadable-kernel-module-lkm-5eaad4db346a1>

³³⁰ <https://man7.org/linux/man-pages/man8/insmod.8.html>

³³¹ https://linux.die.net/man/2/finit_module

³³² <https://github.com/kmod-project/kmod/blob/master/tools/insmod.c>

³³³ <https://www.computerhope.com/unix/insmod.htm>

³³⁴ <https://medium.com/@boutnaru/the-linux-process-journey-dmesg-print-control-the-kernel-ring-buffer-dc78abeb87b7>

³³⁵ <https://slideplayer.com/slide/15437239/>

rmmmod (Remove Kernel Module)

“rmmmod” is a symbolic link³³⁶ to that is located at “/sbin/rmmmod” or “/usr/sbin/rmmmod”. It references the ELF binary “kmod” which is located at “/bin/kmod” or “/usr/bin/kmod”³³⁷. rmmod is used for removing a LKM³³⁸ into the Linux kernel. We can use the “-f” (“--force”) switch of “rmmod” for removing modules which are being used\which are not designed to be removed\have been marked as unsafe. This is relevant only if the kernel was compiled with “CONFIG_MODULE_FORCE_UNLOAD”³³⁹. Also, the “-s” switch (“--syslog”) is used for sending errors to syslog instead of standard error³⁴⁰.

Overall, in order to unload\remove a kernel module we can use the “delete_module” syscall. After glibc version 2.23 there is no wrapper for “delete_module”³⁴¹. Hence, we need to use the “syscall” for indirect syscall invocation³⁴². I suggest going over the code for more details³⁴³. Using rmmod causes the module’s exit function to be called before unloading the LKM³⁴⁴ - as shown in the screenshot below³⁴⁵.

Lastly, until kernel 2.4 the “delete_module” syscall³⁴⁶ took only one argument³⁴⁷. As of the newer kernel it takes two arguments³⁴⁸. By the way, the syscall’s implementation calls “free_module”³⁴⁹ after performing different checks and operations. For more information I suggest going over the “free_module” implementation³⁵⁰.

```
debian@npi:/mnt$ ls
platform_device.ko  platform_driver.ko
debian@npi:/mnt$ sudo insmod platform_device.ko
[ 1360.553491] led_device_init
debian@npi:/mnt$ sudo insmod platform_driver.ko
[ 1367.615716] led_driver_init
[ 1367.619280] led_probe
[ 1367.621822] led major:244
debian@npi:/mnt$ ls /dev/my_led
/dev/my_led
debian@npi:/mnt$ sudo sh -c "echo on > /dev/my_led"
debian@npi:/mnt$ sudo sh -c "echo off > /dev/my_led"
debian@npi:/mnt$ sudo rmmod platform_driver.ko
debian@npi:/mnt$ sudo rmmod platform_device.ko
[ 1439.791876] led device release!
debian@npi:/mnt$ ls /dev/my_led
ls: cannot access '/dev/my_led': No such file or directory
```

CSDN @Couvrire洪荒猛兽

³³⁶ <https://medium.com/@boutnaru/the-linux-concept-journey-linux-file-types-4cb622887331>

³³⁷ <https://medium.com/@boutnaru/the-linux-process-journey-kmod-linux-kernel-module-handling-8ec844216436>

³³⁸ <https://medium.com/@boutnaru/the-linux-concept-journey-loadable-kernel-module-lkm-5eaad4db346a1>

³³⁹ https://cateee.net/lkddb/web-lkddb/MODULE_FORCE_UNLOAD.html

³⁴⁰ <https://man7.org/linux/man-pages/man8/insmod.8.html>

³⁴¹ https://man7.org/linux/man-pages/man2/delete_module.2.html

³⁴² <https://man7.org/linux/man-pages/man2/syscall.2.html>

³⁴³ <https://github.com/kmod-project/kmod/blob/master/tools/rmmod.c>

³⁴⁴ <https://slideplayer.com/slide/15437239/>

³⁴⁵ <http://www.wiks.cn/news/2246.html>

³⁴⁶ https://linux.die.net/man/2/delete_module

³⁴⁷ <https://elixir.bootlin.com/linux/2.4.0/source/include/asm-arm/unistd.h#L407>

³⁴⁸ <https://elixir.bootlin.com/linux/v6.16/source/kernel/module/main.c#L748>

³⁴⁹ <https://elixir.bootlin.com/linux/v6.16/source/kernel/module/main.c#L817>

³⁵⁰ <https://elixir.bootlin.com/linux/v6.16/source/kernel/module/main.c#L1367>

depmod (Dependency Modules)

“depmod” (Dependency Modules) is a symbolic link³⁵¹ to that is located at “/sbin/depmod” or “/usr/sbin/depmod”. It references the ELF binary “kmod” which is located at “/bin/kmod” or “/usr/bin/kmod”³⁵². “depmod” is used for generating a list of dependency descriptions of kernel modules and its associated map files. This is done by going over the “/lib/modules/{kernel-release}” and creating the “modules.dep” file based on the symbols present in the set of modules³⁵³ - as shown in the output of the strace utility below.

Overall, “depmod” is also leveraged for building mappings between hardware identifiers and the modules which handle them. Hence, it can be used for finding the correct module to load when a specific hardware is identified³⁵⁴. “depmod” has different switches like (but not limited to): probing all modules and building the dependency map (“-a”/”-all”), changing input/output directories, report unresolved symbols (“-e”/”-erryms”) and report on symbol version mismatches (“-E Module.symvers”/”--symvers=Module.symvers”). For more information I suggest going over “depmod” documentation³⁵⁵.

Lastly, it is important to understand that beside creating “modules.dep” also a binary hashed version file named “modules.dep.bin” is created in the same directory³⁵⁶. For better understanding I suggest going over the source code of “depmod” as part of kmod’s Github repository³⁵⁷.

```
File Edit View Search Terminal Help
openat(AT_FDCWD, "/lib/modules/[REDACTED]-generic", O_RDONLY|O_NONBLOCK|O_CLOEXEC|O_DIRECTORY) = 3
openat(3, "kernel", O_RDONLY) = 4
openat(4, "net", O_RDONLY) = 5
openat(5, "smc", O_RDONLY) = 6
openat(5, "l2tp", O_RDONLY) = 6
openat(5, "rds", O_RDONLY) = 6
openat(5, "ipv4", O_RDONLY) = 6
openat(6, "netfilter", O_RDONLY) = 7
openat(5, "x25", O_RDONLY) = 6
openat(5, "can", O_RDONLY) = 6
openat(6, "j1939", O_RDONLY) = 7
openat(5, "mac802154", O_RDONLY) = 6
openat(5, "tipc", O_RDONLY) = 6
openat(5, "xdp", O_RDONLY) = 6
openat(5, "ipv6", O_RDONLY) = 6
openat(6, "netfilter", O_RDONLY) = 7
openat(6, "ila", O_RDONLY) = 7
openat(5, "802", O_RDONLY) = 6
openat(5, "bridge", O_RDONLY) = 6
openat(6, "netfilter", O_RDONLY) = 7
openat(5, "sunrpc", O_RDONLY) = 6
openat(6, "xprtrdma", O_RDONLY) = 7
:
```

³⁵¹ <https://medium.com/@boutnaru/the-linux-concept-journey-linux-file-types-4cb622887331>

³⁵² <https://medium.com/@boutnaru/the-linux-process-journey-kmod-linux-kernel-module-handling-8ec844216436>

³⁵³ <https://www.geeksforgeeks.org/linux-unix/depmod-command-in-linux-with-examples/>

³⁵⁴ <https://stackoverflow.com/questions/8697132/why-is-depmod-necessary-for-building-and-working-with-kernel-modules>

³⁵⁵ <https://man.archlinux.org/man/depmod.8.en>

³⁵⁶ <https://man7.org/linux/man-pages/man8/depmod.8.html>

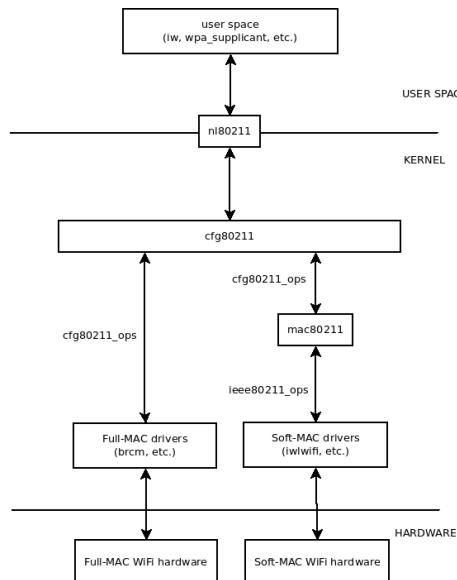
³⁵⁷ <https://github.com/kmod-project/kmod/blob/master/tools/depmod.c>

cfg80211 (Wireless Configuration)

“cfg80211” is a Linux kernel thread which is based on an ordered workqueue³⁵⁸. It is defined as part of “/net/wireless/core.c” in the “cfg80211_init()” function, it is the part of the Linux wireless configuration interface³⁵⁹. It interfaces with “nl80211”, together they are basically a replacement of wext aka” Wireless Extensions”³⁶⁰.

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 stack³⁶¹. “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 interfaces³⁶².

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 scanning³⁶³. Lastly, “cfg80211” interfaces with “nl80211” (as shown in the below) using netlink and not ioctl³⁶⁴. By the way, the name “nl80211” represents 802.11 netlink interface³⁶⁵.



³⁵⁸ <https://elixir.bootlin.com/linux/v6.7.5/source/net/wireless/core.c#L1719>

³⁵⁹ <https://elixir.bootlin.com/linux/v6.7.5/source/net/wireless/core.c#L14>

³⁶⁰ https://elinux.org/images/7/75/DeRosier_WirelessInterfacing.pdf

³⁶¹ <https://johannes.sipsolutions.net/Projects/>

³⁶² <https://elixir.bootlin.com/linux/v6.7.5/source/include/net/cfg80211.h#L4501>

³⁶³ https://wireless.wiki.kernel.org/_media/en/developers/documentation/control.pdf

³⁶⁴ <https://stackoverflow.com/questions/21456235/how-do-the-nl80211-library-cfg80211-work>

³⁶⁵ <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 workqueue³⁶⁶. 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 problematic³⁶⁷.

It is part of the “Device Mapper” framework and used in order to process deferred work hat it has queued up from other contexts where doing immediately so would be problematic³⁶⁸. 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 disk³⁶⁹.

Lastly, there are also other kernel parameters which affect the handling of “dirty” pages, we can check them out using the “sysctl” utility³⁷⁰ - 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 minor³⁷¹.

```
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:~# _
```

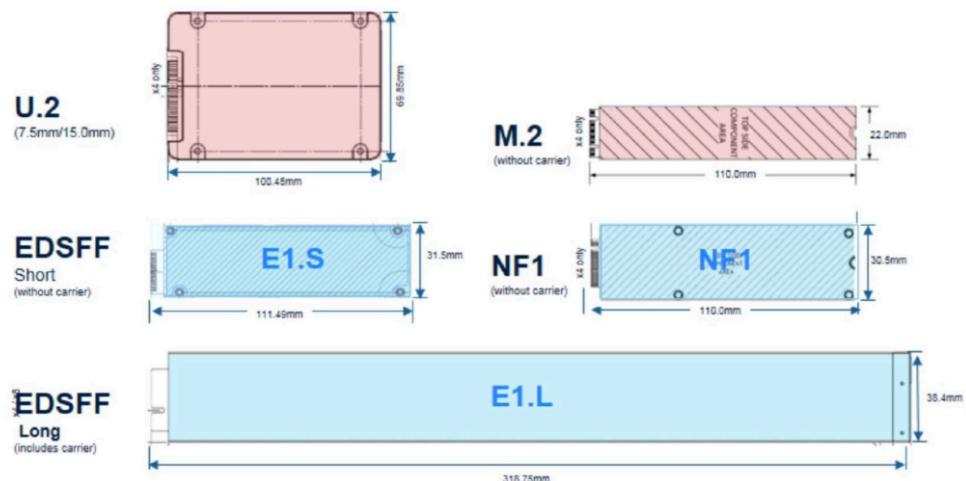
³⁶⁶ <https://elixir.bootlin.com/linux/v6.5-rc3/source/drivers/md/dm.c#L2131>
³⁶⁷ <https://www.compuhov.com/how-do-you-check-which-process-is-using-more-disk-in-linux/>
³⁶⁸ <https://askubuntu.com/questions/986211/what-is-kdmflush>
³⁶⁹ <https://www.cnblogs.com/cobbliu/articles/11792193.html>
³⁷⁰ <https://linux.die.net/man/8/sysctl>
³⁷¹ <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”³⁷². It is used for hosting NVMe related work which is not reset or delete³⁷³. 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 below³⁷⁴. By the way, NVM Express is a non-profit consortium³⁷⁵.

Lastly, “nvme-wq” is used to host work such as: scanning, AEN (Asynchronous Event Notifications) handling, FW activation, periodic reconnects and keep-alive³⁷⁶. Examples of those operations being queued from the Linux source are the following functions: “nvme_queue_scan”³⁷⁷, “nvme_enable_aen”³⁷⁸, and “nvme_fw_act_work”³⁷⁹.



³⁷² <https://elixir.bootlin.com/linux/v6.8.6/source/drivers/nvme/host/core.c#L4831>

³⁷³ <https://elixir.bootlin.com/linux/v6.8.6/source/drivers/nvme/host/core.c#L93>

³⁷⁴ <https://venturebeat.com/business/meet-edssff-1pb-of-flash-storage-in-a-single-rack/>

³⁷⁵ <https://nvmexpress.org/>

³⁷⁶ <https://elixir.bootlin.com/linux/v6.8.6/source/drivers/nvme/host/core.c#L97>

³⁷⁷ <https://elixir.bootlin.com/linux/v6.8.6/source/drivers/nvme/host/core.c#L136>

³⁷⁸ <https://elixir.bootlin.com/linux/v6.8.6/source/drivers/nvme/host/core.c#L1682>

³⁷⁹ <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” utility³⁸⁰. Thus, it is used for checking file types and comparing values³⁸¹. The binary is part of the “coreutils” package³⁸² - 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 file³⁸³, check if a file has been modified since it was last read, if the file exists and the user has read/execute access and more³⁸⁴. 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 command³⁸⁵ - 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 results³⁸⁶.

```
root@localhost:~# builtin [
-bash: [: missing `]'
root@localhost:~# builtin /usr/bin/[[
-bash: builtin: /usr/bin/[ is not a shell builtin
root@localhost:~# file /usr/bin/[[
/usr/bin/[ ELF 32-bit LSB pie executable, Intel 80386, version 1 (SYSU), 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 /dcfugbhnmds ]
root@localhost:~# echo $?
1
```

³⁸⁰ <https://superuser.com/questions/334549/what-is-usr-bin-and-how-do-i-use-it>

³⁸¹ <https://www.man7.org/linux/man-pages/man1/test.1.html>

³⁸² <https://github.com/coreutils/coreutils/blob/master/src/test.c#L36>

³⁸³ <https://medium.com/@boutnaru/the-linux-concept-journey-linux-file-types-4cb622887331>

³⁸⁴ <https://linux.die.net/man/1/test>

³⁸⁵ <https://linux.die.net/man/1/builtins>

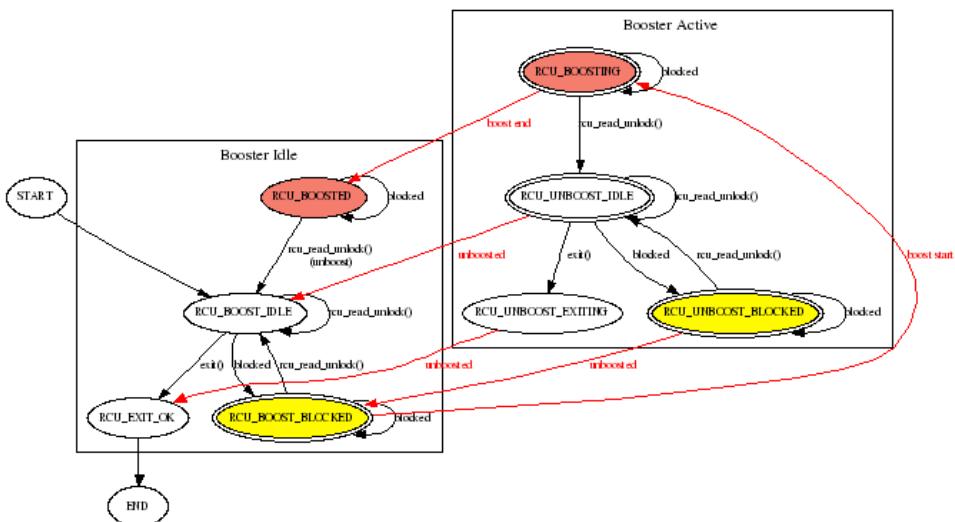
³⁸⁶ [https://grep.app/search?q=%5B%20&filter\[path\]\[0\]=scripts/](https://grep.app/search?q=%5B%20&filter[path][0]=scripts/)

rcub (Read-Copy Update Boost)

“rcub” is a Linux kernel thread created as part of the “rcu_spawn_one_boost_kthread” function³⁸⁷. 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 RCU³⁸⁸.

Overall, RCU (Read-Copy Update) is a mutual exclusion mechanism³⁸⁹. 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 below³⁹⁰. We can enable RCU priority boosting using the “CONFIG_RCU_BOOST” configuration item when compiling the kernel³⁹¹.

Lastly, The “rcub” kernel thread executes the code of the “rcu_boost_kthread” function³⁹². 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 loads³⁹³.



³⁸⁷ https://elixir.bootlin.com/linux/v6.8.5/source/kernel/rcu/tree_plugin.h#L1202

³⁸⁸ https://elixir.bootlin.com/linux/v6.8.5/source/kernel/rcu/tree_plugin.h#L1188

³⁸⁹ <https://medium.com/@boutnaru/linux-rCU-read-copy-update-c2793ee7a4cd>

³⁹⁰ <https://lwn.net/Articles/220677/>

³⁹¹ <https://lwn.net/Articles/777214/>

³⁹² https://elixir.bootlin.com/linux/v6.8.5/source/kernel/rcu/tree_plugin.h#L1110

³⁹³ https://cateee.net/lkdb/web-lkdb/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 below³⁹⁴.

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 kernel³⁹⁵.

Moreover, “dmesg” is based on reading information for the character device “/dev/kmsg”³⁹⁶. 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 more³⁹⁷.

Lastly, “dmesg” is part of the “util-linux” package, we can go over the source code of the utility in repo of the package³⁹⁸. We can also force “dmesg” to use the “syslog” syscall³⁹⁹ to read kernel messages instead of using “/dev/kmsg”⁴⁰⁰. 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=
```

³⁹⁴ <https://linuxize.com/post/dmesg-command-in-linux/>

³⁹⁵ https://linuxsecurity.expert/kb/sysctl/kernel_dmesg_restrict/

³⁹⁶ <https://www.kernel.org/doc/Documentation/ABI/testing/dev-kmsg>

³⁹⁷ <https://man7.org/linux/man-pages/man1/dmesg.1.html>

³⁹⁸ <https://github.com/util-linux/util-linux/blob/master/sys-utils/dmesg.c>

³⁹⁹ <https://man7.org/linux/man-pages/man2/syslog.2.html>

⁴⁰⁰ <https://github.com/util-linux/util-linux/blob/master/sys-utils/dmesg.c#L343>

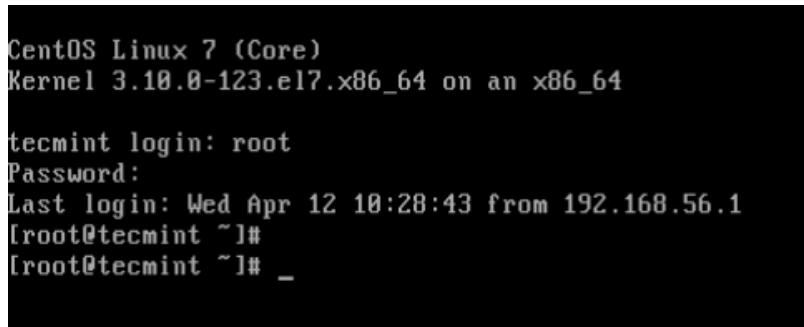
login (Begin Session on The System)

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 shown⁴⁰¹ - as shown in the screenshot below⁴⁰².

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 more⁴⁰³. We can check the source code of “login” as part of the “util-linx” github repository⁴⁰⁴.

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”⁴⁰⁵ which will show libpam.so and probably also libpam_misc.so⁴⁰⁶.

Lastly, there are different configuration files that affect the behavior of “loing” (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 writeups⁴⁰⁷.



The screenshot shows a terminal window on a CentOS Linux 7 system. The title bar indicates "CentOS Linux 7 (Core)" and "Kernel 3.10.0-123.el7.x86_64 on an x86_64". The terminal displays a root login session:

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

⁴⁰¹ https://man7.org/linux/man-pages/man1/login_1.html

⁴⁰² <https://www.tecmint.com/understanding-shell-initialization-files-and-user-profiles-linux/>

⁴⁰³ <https://en.wikipedia.org/wiki/Util-linux>

⁴⁰⁴ <https://github.com/util-linux/util-linux/blob/master/login-utils/login.c>

⁴⁰⁵ <https://medium.com/@boutnaru/linux-instrumentation-part-4-ldd-888502965a9b>

⁴⁰⁶ <https://medium.com/@boutnaru/the-linux-security-journey-pam-pluggable-authentication-module-388496a8785c>

⁴⁰⁷ <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”⁴⁰⁸. “su” is part of the “util-linux” package, we can go over the implementation of the command as part of the package’s Github repository⁴⁰⁹.

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 root⁴¹⁰ - 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 variable⁴¹¹. By the way, it is recommended to use “su --login” instead of “su -” as a shortcut to avoid side effects caused by mixing environments⁴¹².

```
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:
```

⁴⁰⁸ <https://linuxize.com/post/su-command-in-linux/>

⁴⁰⁹ <https://github.com/util-linux/util-linux/blob/master/login-utils/su-common.c>

⁴¹⁰ <https://www.tecmint.com/difference-between-su-and-su-commands-in-linux/>

⁴¹¹ <https://www.geeksforgeeks.org/difference-between-su-and-su-command-in-linux/>

⁴¹² <https://man7.org/linux/man-pages/man1/su.1.html>

ps (Process Status)

“ps” is an ELF binary located at “/usr/bin/ps” (or “/bin/ps”) it is used for retrieving a snapshot of the current processes\threads\tasks⁴¹³. By the way, the command is relevant for various operating systems (not only Linux) such as: “Unix”, “Plan 9”, “IBM i”, “KolibriOS” and more⁴¹⁴.

Overall, “ps” is based on reading the information it reports from “/proc”. Also, in case we want to override the default output format of the command by leveraging the “PS_FORMAT” environment variable, which is not the only one that affects “ps”⁴¹⁵. “ps” on Linux supports both a standard syntax and a BSD syntax. For example “ps -e”\”ps -ef”\”ps -ely” shows all processes on the system in standard syntax while “ps ax”\”ps aux” does the same using the BSD syntax⁴¹⁶ - as shown in the syntax below (taken using <https://copy.sh/v86/?profile=archlinux>).

Lastly, among the information that “ps” can displayed we can find the following: CPU usage, memory usage, nice value, pid, ppid, priority, status code, time when the process started, associated terminal, uid and more⁴¹⁷. “ps” is also one of the utilities which is included as part of “busybox”⁴¹⁸. For a reference implementation we can check out the “procps” repository in GitLab⁴¹⁹.

```
root@localhost:~# ps aux | wc -l
58
root@localhost:~# ps -ef | wc -l
58
root@localhost:~# ps aux | head -5
USER      PID %CPU %MEM    VSZ   RSS TTY      STAT START   TIME COMMAND
root        1  0.3  0.3  2596  1580 ?        Ss  12:19  0:01 init [3]
root        2  0.0  0.0     0   0 ?        S   12:19  0:00 [kthreadd]
root        3  0.0  0.0     0   0 ?        I<  12:19  0:00 [rcu_gp]
root        4  0.0  0.0     0   0 ?        I<  12:19  0:00 [rcu_par_gp]
root@localhost:~# ps -ef | head -5
UID      PID  PPID  C STIME TTY          TIME CMD
root        1    0  0 12:19 ?        00:00:01 init [3]
root        2    0  0 12:19 ?        00:00:00 [kthreadd]
root        3    2  0 12:19 ?        00:00:00 [rcu_gp]
root        4    2  0 12:19 ?        00:00:00 [rcu_par_gp]
root@localhost:~# _
```

⁴¹³ <https://man7.org/linux/man-pages/man1/ps.1.html>

⁴¹⁴ [https://en.wikipedia.org/wiki/Ps_\(Unix\)](https://en.wikipedia.org/wiki/Ps_(Unix))

⁴¹⁵ <https://wiki.gentoo.org/wiki/Ps>

⁴¹⁶ <https://man.archlinux.org/man/ps.1.en>

⁴¹⁷ [https://www.wikiwand.com/en/articles/Ps_\(Unix\)](https://www.wikiwand.com/en/articles/Ps_(Unix))

⁴¹⁸ <https://medium.com/@boutnaru/the-linux-concept-journey-busybox-234cfdc8808a>

⁴¹⁹ <https://gitlab.com/procps-ng/procps>

kacpid (Kernel Advanced Configuration and Power Interface Daemon)

“kacpid” is a Linux kernel thread which is based on an workqueue⁴²⁰. It is defined as part of “drivers/acpi/osl.c” in the “acpi_os_initialize1()” function, as part of the ACPI support of Linux⁴²¹.

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 below⁴²². Also, it can be used for checking battery levels, PCI IRQ routing, CPUs, NUMA domains, fan speeds, temperature sensors and more⁴²³.

Lastly, work is queued for the “kacpid” kernel as part of the “acpi_os_execute()” function⁴²⁴. This kernel thread is used in case of a GPE (General Purpose Event) handler is needed⁴²⁵. 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 laptop⁴²⁶.

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.
	(S1)	Requires less power than the G0 state and has multiple sleeping states.
	(S2)	CPU is still powered, and unused devices are powered down. RAM is still being refreshed. Hard disks are not running.
	(S3)	CPU is not powered. RAM is still being refreshed. System is restored instantly upon user intervention.
	(S4)	Power supply output is reduced. RAM is still being refreshed. Some info in RAM is restored to CPU and cache.
G2	(S5)	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.
G3		Also called soft off. Power consumption is almost zero. Requires the operating system to reboot. No information is saved anywhere.
		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.

⁴²⁰ <https://elixir.bootlin.com/linux/v6.9.5/source/drivers/acpi/osl.c#L1665>

⁴²¹ <https://elixir.bootlin.com/linux/v6.9.5/source/drivers/acpi>

⁴²² <https://www.slideserve.com/reecce/chapter-4-disassembly-and-power>

⁴²³ <https://wiki.osdev.org/ACPI>

⁴²⁴ <https://elixir.bootlin.com/linux/v6.9.5/source/drivers/acpi/osl.c#L1113>

⁴²⁵ <https://elixir.bootlin.com/linux/v6.9.5/source/drivers/acpi/osl.c#L1105>

⁴²⁶ <https://askubuntu.com/questions/148726/what-is-an-acpi-gpe-storm>

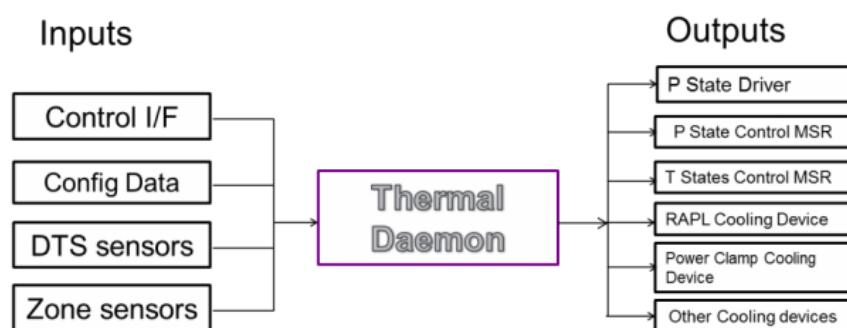
thermald (Thermal Daemon)

“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 value⁴²⁷. We can go over the source code of the thermal daemon in GitHub⁴²⁸.

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 points and based on predefined configuration use the best cooling method⁴²⁹ - 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” CPUs⁴³⁰.

Block Diagram



⁴²⁷ <https://wiki.debian.org/thermald>

⁴²⁸ https://github.com/intel/thermal_daemon

⁴²⁹ <https://web.archive.org/web/20211123224127/https://01.org/linux-thermal-daemon/documentation/introduction-thermal-daemon>

⁴³⁰ <https://www.linux.com/news/linux-thermal-daemon-monitors-and-controls-temperature-tablets-laptops/>

dhcpd (Dynamic Host Configuration Protocol Daemon)

“dhcpd” is an ELF binary located by default at “/usr/bin/dhcpd” (or “/bin/dhcpd”). 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)⁴³¹.

Overall, we can find the “dhcpd” configuration in the “dhcpd.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 more⁴³². “dhcpd” 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”⁴³³. By the way, the ISC replacement is called “Kea” (more on that in a future writeup).

Lastly, it is important to understand that “dhcpd” and “dhpcd” are not the same. The first is a server implementation while the second is a client implementation of DHCP. Also, “dhcpd” includes both an implementation for DHCPv4 and DHCPv6⁴³⁴. More information about “dhcpd” can be found as part of the ISC documentation⁴³⁵. On some Linux distributions we can use “systemctl”⁴³⁶ for performing different operations of the “dhcp” daemon like checking its status - as shown in the screenshot below⁴³⁷.

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

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

Ishovom@linuxhint-dhcp-server-6c024 ~]$
```

⁴³¹ <https://www.spiceworks.com/tech/networking/articles/what-is-dhcp/>

⁴³² <https://linux.die.net/man/5/dhcpd.conf>

⁴³³ <https://www.isc.org/blogs/isc-dhcp-eol/>

⁴³⁴ <https://wiki.archlinux.org/title/dhcpd>

⁴³⁵ <https://kb.isc.org/docs/aa-00333>

⁴³⁶ <https://www.man7.org/linux/man-pages/man1/systemctl.1.html>

⁴³⁷ https://linuxhint.com/dhcp_server_centos8/

ata_sff (Advanced Technology Attachment Small Form Factor)

“ata_sff” is a Linux kernel thread which is based on a workqueue⁴³⁸ - as shown in the screenshot below (taken from <https://copy.sh/v86/?profile=archlinux>). The source code file continuing that (libata-sff.c) states it is used as a helper library for PCI IDE BMDMA⁴³⁹. By the way, BMDMA stands for “Bus-Master Direct Memory Access”⁴⁴⁰.

Overall, to enable “ata_sff” support (for legacy IDE and PATA) we need to enable “CONFIG_ATA_SFF” while compiling the kernel⁴⁴¹. Also, ATA originally was designed to work with hard disks (and devices that could emulate them). However a committee called SFF (Small Form Factor) (SFF) introduced the ATAPI which can be used for a variety of other devices, which includes functions beyond those necessary for hard disks⁴⁴².

Moreover, libATA is a library used inside the Linux kernel to support ATA (Advanced Technology Attachment) host controllers/devices. It is used for providing: class transport for ATA/ATAPI devices, SCSI<->ATA translation based on T10 SAT specification for ATA devices and ATA driver API⁴⁴³. Among the features that are supported are: power management, analysis and reporting, hot swapping, NCQ (native command queueing) and port multiplying⁴⁴⁴. I suggest going over the libATA’s developer guide for more information about the API⁴⁴⁵.

Lastly, work is queued to “ata_sff” using the “ata_sff_queue_work” function⁴⁴⁶. An example of using it can be found as of the Arasan Compact Flash host controller source⁴⁴⁷. By the way, there is also a possibility to queue the work after a specific delay using the “ata_sff_queue_delayed_work” function⁴⁴⁸.

```
root@localhost:~# ps -ef | grep -v grep | grep ata_sff
root      33      2  0 02:50 ?          00:00:00 [ata_sff]
root@localhost:~# cat /proc/33/stack
[<0>] rescuer_thread+0x273/0x350
[<0>] kthread+0xcc/0xf0
[<0>] ret_from_fork+0x19/0x28
```

⁴³⁸ <https://elixir.bootlin.com/linux/v6.7.4/source/drivers/ata/libata-sff.c#L3197>

⁴³⁹ <https://elixir.bootlin.com/linux/v6.7.4/source/drivers/ata/libata-sff.c#L3>

⁴⁴⁰ <https://www.linux.org/threads/the-linux-kernel-configuring-the-kernel-part-13.9143/>

⁴⁴¹ https://cateee.net/lkddb/web-lkddb/ATA_SFF.html

442 <https://en.wikipedia.org/wiki/ATAPI>443 <https://www.kernel.org/doc/Documentation/driver-api/libata.rst>444 <https://www.linux-ata.org/features.html>445 <https://www.kernel.org/doc/html/v4.14/driver-api/libata.html>446 <https://elixir.bootlin.com/linux/v6.7.4/source/drivers/ata/libata-sff.c#L1165>447 https://elixir.bootlin.com/linux/v6.7.4/source/drivers/ata/pata_arasan_cf.c#L686448 <https://elixir.bootlin.com/linux/v6.7.4/source/drivers/ata/libata-sff.c#L1171>

uptime (System's Uptime)

“uptime” is a Linux ELF binary located “/bin/uptime” and/or “/usr/bin/uptime”. It is used for telling how long the specific Linux system has been running. This means the amount of time a system remains operational and accessible without interruptions⁴⁴⁹. Among the information which is displayed by uptime we can find the following: the current time, how much the system is running, how many users are logged on and the system load averages for the past 5/10/15 minutes⁴⁵⁰ - as shown in the screenshot below (taken from copy.sh).

Overall, for collecting information about the logged on user “uptime” accesses “/var/run/utmp”. In case of the system load averages “/proc/loadavg”⁴⁵¹ file is used for getting the information. To get the total time the system is running the “/proc/uptime”⁴⁵² file is accessed - as shown in the output of the system call tracing (using strace) below.

Lastly, “uptime” also has command line switches that can be used. Examples are: “-p” which prints the output of the uptime in human readable format and “-s” which prints the full timestamp from which the system is alive⁴⁵³ - as shown in the screenshot below. “uptime” is part of the “coreutils” package, we can also check the source code for more information⁴⁵⁴.

```
root@localhost:~# uptime
04:13:21 up 1:23, 2 users,  load average: 0.00, 0.00, 0.00
root@localhost:~# uptime -s
2022-11-07 02:50:13
root@localhost:~# uptime -p
up 1 hour, 23 minutes
root@localhost:~# strace -e openat
openat(AT_FDCWD, "/etc/ld.so.cache", O_RDONLY|O_LARGEFILE|O_CLOEXEC) = 3
openat(AT_FDCWD, "/usr/lib/libprocps.so.0", O_RDONLY|O_LARGEFILE|O_CLOEXEC) = 3
openat(AT_FDCWD, "/usr/lib/libc.so.6", O_RDONLY|O_LARGEFILE|O_CLOEXEC) = 3
openat(AT_FDCWD, "/usr/lib/libcrypt.so.0", O_RDONLY|O_LARGEFILE|O_CLOEXEC) = 3
openat(AT_FDCWD, "/usr/lib/libdl.so.2", O_RDONLY|O_LARGEFILE|O_CLOEXEC) = 3
openat(AT_FDCWD, "/usr/lib/libcap.so.2", O_RDONLY|O_LARGEFILE|O_CLOEXEC) = 3
openat(AT_FDCWD, "/usr/lib/libcrypt.so.20", O_RDONLY|O_LARGEFILE|O_CLOEXEC) = 3
openat(AT_FDCWD, "/usr/lib/liblzo.so.5", O_RDONLY|O_LARGEFILE|O_CLOEXEC) = 3
openat(AT_FDCWD, "/usr/lib/libzstd.so.1", O_RDONLY|O_LARGEFILE|O_CLOEXEC) = 3
openat(AT_FDCWD, "/usr/lib/liblz4.so.1", O_RDONLY|O_LARGEFILE|O_CLOEXEC) = 3
openat(AT_FDCWD, "/usr/lib/libgcc_s.so.1", O_RDONLY|O_LARGEFILE|O_CLOEXEC) = 3
openat(AT_FDCWD, "/usr/lib/libgpg-error.so.0", O_RDONLY|O_LARGEFILE|O_CLOEXEC) = 3
openat(AT_FDCWD, "/proc/self/auxv", O_RDONLY|O_LARGEFILE) = 3
openat(AT_FDCWD, "/proc/sys/kernel/osrelease", O_RDONLY|O_LARGEFILE) = 3
openat(AT_FDCWD, "/sys/devices/system/cpu/online", O_RDONLY|O_CLOEXEC) = 3
openat(AT_FDCWD, "/proc/self/auxv", O_RDONLY|O_LARGEFILE) = 3
openat(AT_FDCWD, "/etc/localtime", O_RDONLY|O_CLOEXEC) = -1 ENOENT (No such file or directory)
openat(AT_FDCWD, "/proc/uptime", O_RDONLY|O_LARGEFILE) = 3
openat(AT_FDCWD, "/var/run/utmp", O_RDONLY|O_LARGEFILE|O_CLOEXEC) = 4
openat(AT_FDCWD, "/proc/loadavg", O_RDONLY|O_LARGEFILE) = 4
04:13:45 up 1:23, 2 users,  load average: 0.00, 0.00, 0.00
*** exited with 0 ***
root@localhost:~#
```

⁴⁴⁹ <https://phoenixnap.com/glossary/what-is-uptime>

⁴⁵⁰ <https://linux.die.net/man/1/uptime>

⁴⁵¹ https://man7.org/linux/man-pages/man5/proc_loadavg.5.html

⁴⁵² https://man7.org/linux/man-pages/man5/proc_uptime.5.html

⁴⁵³ <https://phoenixnap.com/kb/linux-uptime>

⁴⁵⁴ <https://github.com/coreutils/coreutils/blob/master/src/uptime.c>

agetty (Alternative Linux getty)

“getty” is a Linux ELF binary located at “/bin/agetty” and/or “/usr/bin/agetty”. The binary is used for opening a new tty port and prompting for a login by leveraging “/bin/login”⁴⁵⁵. Also, “agetty” is normally launched by PID 0 (init/systemd). Moreover, “agetty” has extra features which can be usual for dial up\hardwired connections such as: deducing baud rate, not displaying the content of /etc/issue, invoking a different login program (not /bin/login), forcing the line to be local with no need for carrier detect and more⁴⁵⁶.

Overall, “agetty” is the Linux version of “getty” (which is short for “get TTY”). In the realm of Unix “getty” is programming running on a host computer which is used for managing physical/virtual terminals for providing multi-user access⁴⁵⁷. We can pass different arguments to “agetty” to configure its operation such as “-8” which is assume 8-bit tty mode - as shown in the screenshot below⁴⁵⁸

Lastly, “agetty” is part of the “util-linux” set of utilities. We can think about it as a replacement for getty on SunOS 4.1.x or the SAC ttymon/ttyadm/sacadm/pmadm suite on Solaris and other SVR4 systems. We can find the code of “agetty” as part of the “util-linux” GitHub repository⁴⁵⁹.

```
ubuntu@dev-vm:~$ agetty -8 - linux
Ubuntu 20.04.4 LTS dev-vm -
dev-vm login: █
```

⁴⁵⁵ <https://medium.com/@boutnaru/the-linux-process-journey-login-02b6d83ab6c5>

⁴⁵⁶ <https://man7.org/linux/man-pages/man8/agetty.8.html>

⁴⁵⁷ <https://www.geeksforgeeks.org/agetty-command-in-linux-with-examples/>

⁴⁵⁸ <https://linuxgenie.net/agetty-command-in-linux/>

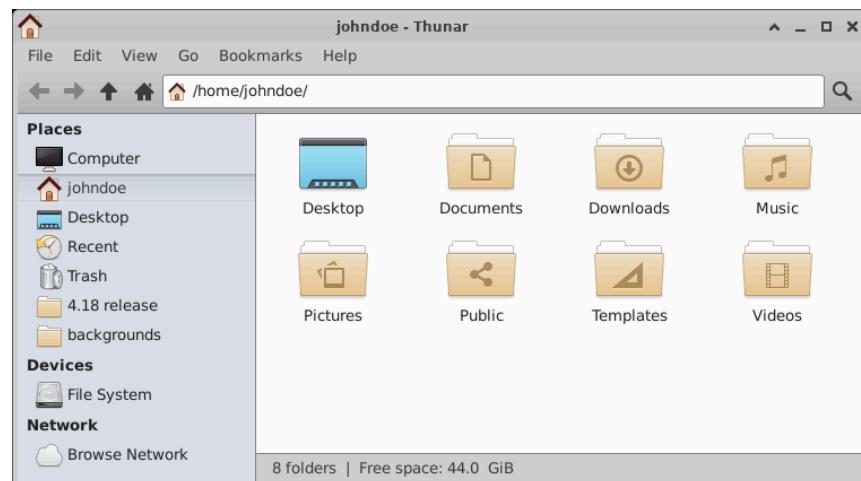
⁴⁵⁹ <https://github.com/util-linux/util-linux/blob/master/term-utils/agetty.c>

Thunar (Xfce File Manager)

“Thunar” is a symbolic link located at “/bin/Thunar” and/or “/usr/bin/Thunar”. It points to “/bin/thunar” and “/usr/bin/thunar” respectively, which are Linux ELF binaries. Thunar provides basic file manipulation and handling like: renaming, deleting, watching attributes, searching and more. Also, Thunar allows working with removal media and network storage such as NFS/SMB⁴⁶⁰.

Overall, Thunar is a modern file manager as part of the Xfce desktop environment⁴⁶¹. We can think about it as similar to Windows’ “File Explorer”⁴⁶² - as shown on the screenshot below⁴⁶³. Thunar also supports installing plugins which are basically additional packages that are used in order to extend functionality⁴⁶⁴.

Lastly, it is highly recommended to go over the source code of thunar⁴⁶⁵. Thunar can be configured to run commands automatically when media are connected (If both gvfs and thunar-volman are installed). In case of mobile devices (which usually are based on MTP) , an additional gvfs-mtp package is required. For Thunar to support automatic mounting, we need to launch thunar in daemon mode. This is done by providing the “--daemon” argument to the thunar executable⁴⁶⁶.



⁴⁶⁰ <https://docs.xfce.org/xfce/thunar/start>

⁴⁶¹ <https://medium.com/@boutnaru/the-linux-concept-journey-xfce-desktop-environment-f91b5272f5ed>

⁴⁶² <https://medium.com/@boutnaru/the-windows-concept-journey-file-explorer-previously-windows-explorer-e48077b135a0>

⁴⁶³ <http://docs.xfce.org/xfce/thunar/the-file-manager-window>

⁴⁶⁴ <https://goodies.xfce.org/projects/thunar-plugins/start>

⁴⁶⁵ <https://gitlab.xfce.org/xfce/thunar>

⁴⁶⁶ <https://wiki.archlinux.org/title/Thunar>

ssh (Secure Shell Client)

“ssh” (Secure Shell Client) is a Linux ELF binary located at “/bin/ssh” and/or “/usr/bin/ssh” (could be also “/bin/ssh” or “/usr/sbin/ssh”). It provides the ability to remote login to a Linux system. After logging in we can execute commands on the remote machine - as shown in the screenshot below⁴⁶⁷. Also, it is used as a replacement for rlogin/rsh. SSH creates a secure encrypted channel between the client and server. By the way, TCP communications (such as X11) can be forwarded over the SSH tunnel⁴⁶⁸.

Overall, the SSH client can authenticate using passwords (less secure) or SSH keys (which is considered to be more secure). When authenticating with SSH keys we need to have the key pair on the local machine. Also, the remote system needs the public key to be stored in “~/.ssh/authorized_keys” (where “~” stands for the home folder of the connecting user). The “authorized_keys” file contains a list of public keys (one public key per line) which are authorized to log into the specific user account⁴⁶⁹.

Lastly, “ssh” also has its own configuration which is stored in “~.ssh/config”⁴⁷⁰. Among the configuration options we can find: enabling X11 forwarding and agent forwarding, port forwarding, configuring public key authentication and more⁴⁷¹. For more information we can check out the OpenSSH GitHub repository⁴⁷².

```
[root@Client ~]# ssh 192.168.126.143
root@192.168.126.143's password:
Last login: Mon Dec  3 12:03:18 2018 from 192.168.126.1
[root@Server ~]# ifconfig
eth0      Link encap:Ethernet HWaddr 00:0C:29:E1:7B:AD
          inet addr:192.168.126.143 Bcast:192.168.126.255 Mask:255.255.255.0
          inet6 addr: fe80::20c:29ff:fe1:7bad/64 Scope:Link
            UP BROADCAST RUNNING MULTICAST MTU:1500 Metric:1
            RX packets:1007 errors:0 dropped:0 overruns:0 frame:0
            TX packets:730 errors:0 dropped:0 overruns:0 carrier:0
            collisions:0 txqueuelen:1000
            RX bytes:86155 (84.1 KiB) TX bytes:78044 (76.2 KiB)

lo       Link encap:Local Loopback
          inet addr:127.0.0.1 Mask:255.0.0.0
          inet6 addr: ::1/128 Scope:Host
            UP LOOPBACK RUNNING MTU:16436 Metric:1
            RX packets:8 errors:0 dropped:0 overruns:0 frame:0
            TX packets:8 errors:0 dropped:0 overruns:0 carrier:0
            collisions:0 txqueuelen:0
            RX bytes:480 (480.0 b) TX bytes:480 (480.0 b)

[root@Server ~]#
```

⁴⁶⁷ <https://www.youtube.com/watch?v=0PRd4aCyrk>

⁴⁶⁸ <https://linux.die.net/man/1/ssh>

⁴⁶⁹ <https://www.digitalocean.com/community/tutorials/ssh-essentials-working-with-ssh-servers-clients-and-keys>

⁴⁷⁰ <https://www.geeksforgeeks.org/how-to-configure-ssh-client-in-linux/>

⁴⁷¹ <https://www.ssh.com/academy/ssh/config>

⁴⁷² <https://github.com/openssh/openssh-portable/blob/master/ssh.c>

sshd (OpenSSH SSH Daemon)

“sshd” (OpenSSH SSH Daemon) is a Linux ELF binary located at “/sbin/sshd” and/or “/usr/sbin/sshd” (In rare cases “/bin/sshd” or “/usr/bin/sshd”). It is the demon part for the “ssh”⁴⁷³ client. “sshd” is the server part which provides a secure encrypted channel between two endpoints over the network. “sshd” is part of the openssh-server Linux package⁴⁷⁴.

Overall, “sshd” listens for incoming connections from clients. By default it listens on port 22 which can be modified by configuration or by using the “-p” argument. Also, it mostly is being run in the background as an operating system service. For every incoming connection the process is forked. The forked instance handles key exchange, encryption, authentication, command execution, and data exchange⁴⁷⁵ - as shown in the screenshot below⁴⁷⁶.

Lastly, the best practice for configuring the SSH server is by leveraging a configuration file which by default is “/etc/ssh/ssh_config”. We can also use a different configuration file by specifying the “-f” argument. Examples for configuration options are: ciphers, chroot directory, allowed users, authorized key files and more⁴⁷⁷. For more information we can check out the OpenSSH GitHub repository⁴⁷⁸.

```
[root@fresh-bytes-1 ~]# pstree
systemd--NetworkManager--dhclient
                                         └─2*[{NetworkManager}]
                                         acpid
                                        agetty
                                        audited--{audited}
                                         chronyd
                                         containerd--10*[{containerd}]
                                         crond
                                         dbus-daemon
                                         dockerd--10*[{dockerd}]
                                         haveged
                                         irqbalance
                                         master--pickup
                                         └─qmgr
                                         polkitd--5*[{polkitd}]
                                         qemu-ga
                                         qemu-kvm_ga--qemu-kvm_ga
                                         rsyslogd--2*[{rsyslogd}]
                                         sshd--sshd--bash--pstree
                                         systemd-journal
                                         └─systemd-logind
                                         └─systemd-udevd
                                         tuned--4*[{tuned}]
[root@fresh-bytes-1 ~]#
```

⁴⁷³ <https://medium.com/@boutnaru/the-linux-process-journey-ssh-secure-shell-client-01cc67897629>

<https://launchpad.net/ubuntu/+source/openssl>

⁴⁷⁵ <https://linux.die.net/man/8/sshd>

476 <https://wumanho.cn/posts/cmd/>

https://linux.die.net/man/5/sshd_config

478 <https://github.com/openssh/openssh-portable>

For more information about the NIST Privacy Framework, visit www.nist.gov/privacy-framework.

passwd (Change User Password)

“passwd” is an ELF binary located at “/usr/bin/passwd” (or “/bin/passwd”). It is used for changing the password for a specific user account. Due to the fact it needs to modify root only writable content (“/etc/shadow” and/or “/etc/passwd”) the binary is owned by root and marked with “suid” bit⁴⁷⁹.

Overall, every user can change only its own password thanks to the “real user ID” aka RUID⁴⁸⁰ - as shown in the screenshot below. Root\supersuer can of course alter the password of any user (given as an argument to the command). The user is first prompted for its old password which is verified, if the password is correct the user can modify the password. In case of a superuser\root the first phase is bypassed to allow resetting a password if the user has forgotten it. Also, the password is tested for its complexity for security reasons⁴⁸¹.

Lastly, when the user changes a password a prompt is presented twice. The second time is compared to the first time in order to ensure the two password input match (this is a safety mechanism). The password is checked against different requirements like: length, history checking, ensuring it is not part of a configured dictionary and more⁴⁸². This can be done using the “/etc/default/passwd” configuration file and\or PAM⁴⁸³.

The screenshot shows a terminal window titled "Terminal Session". The terminal output is as follows:

```
Troller $ ls -lah `which passwd`  
-rwsr-xr-x 1 root root 59K Mar 14 2022 /usr/bin/passwd  
Troller $ id | cut -d "(" -f1  
uid=1000  
Troller $ passwd  
Changing password for [REDACTED]  
Current password: [REDACTED]  
Troller $ ps -a -o comm,euid,ruid,suid,fsuid,egid,rgid,sgid,fsgid  
COMMAND EUID RUID SUID FSUID EGID RGID SGID FSGID  
passwd 0 1000 0 0 1000 1000 1000 1000  
ps 1000 1000 1000 1000 1000 1000 1000 1000  
Troller $
```

⁴⁷⁹ <https://medium.com/@boutnaru/linux-security-suid-bit-d4f553e7d99e>

⁴⁸⁰ <https://medium.com/@boutnaru/the-linux-security-journey-ruid-real-user-id-b23abcbea9c6>

⁴⁸¹ <https://man7.org/linux/man-pages/man1/passwd.1.html>

⁴⁸² https://docs.oracle.com/cd/E88353_01/html/E37839/passwd-1.html

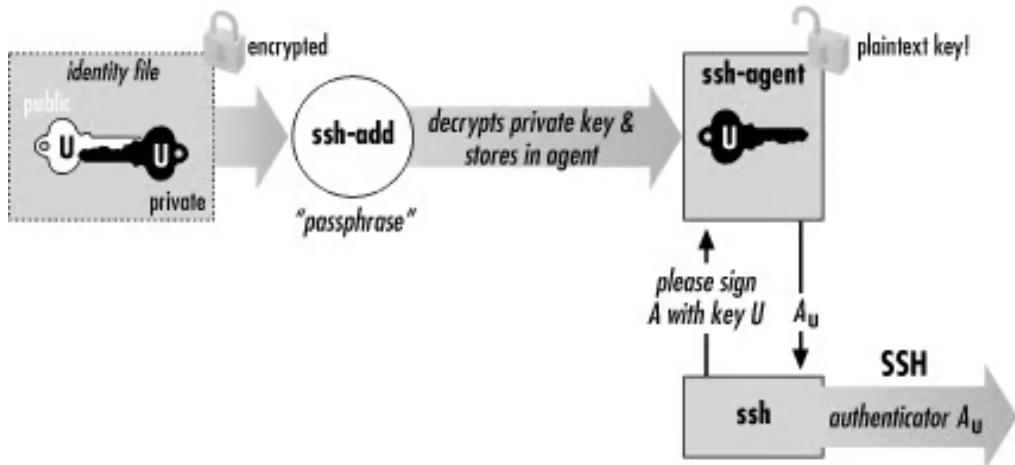
⁴⁸³ <https://medium.com/@boutnaru/the-linux-security-journey-pam-pluggable-authentication-module-388496a8785c>

ssh-agent (Secure Shell Authentication Agent)

“ssh-agent” (Secure Shell Authentication Agent) is a Linux ELF binary located at “/bin/ssh-agent” and/or “/usr/bin/ssh-agent”. It is an integral part of the OpenSSH package. “ssh-agent” is the key manager for SSH. This means it holds keys/certificates in memory in an unencrypted manner ready for ssh to use (saves us retying passphrases every time we connect to a server). It runs separately from the ssh binary in the background after running ssh for the first time⁴⁸⁴ - as shown in the diagram below⁴⁸⁵.

Overall, the reason “ssh-agent” keeps crypto materials (RSA/DSA/etc) safe is due to the fact they are not exported nor written to disk. “ssh-agent” is started at the beginning of an X-session or a login session. Thus, all other windows or programs are started as clients of it. By using environment variables “ssh-agent” can be located and used for authentication when logging using SSH⁴⁸⁶.

Lastly, an ssh client and the “ssh-agent” are “speaking” to each other using the “SSH Agent Protocol”⁴⁸⁷ - more on that in future writeups. The basic operations that “ssh-agent” is providing are: adding a key pair, removing a key/key pair, sign messages using the stored keys, listing keys stored and more. We can check out OpenSSH’s GitHub repository for more information⁴⁸⁸.



⁴⁸⁴ <https://smallstep.com/blog/ssh-agent-explained/>

⁴⁸⁵ <https://dev.to/samuvil/ssh-agents-in-depth-4116>

⁴⁸⁶ <https://linux.die.net/man/1/ssh-agent>

⁴⁸⁷ <https://datatracker.ietf.org/doc/html/draft-miller-ssh-agent-04>

⁴⁸⁸ <https://github.com/openssh/openssh-portable/blob/master/ssh-agent.c>

hwrng (Hardware Number Random Generator)

“hwrng” is a Linux kernel thread which is executed using the “kthread_run” (a wrapper for ktrehad_create) function⁴⁸⁹. The specific function which is used for that is “hwrng_fillfn”⁴⁹⁰. It is provided by the “hw_random” framework, which is a software component that makes use of a hardware based RNG (Random Number Generator) on the CPU\motherboard. It has two major parts: a character device⁴⁹¹ “/dev/hwrng” including a sysfs support and a hardware-specific driver⁴⁹².

Overall, “/dev/hwrng” has the following <major,minor>⁴⁹³ number pair: <10, 183> and its relevant sysfs attributes are located at “/sys/class/misc/hw_random” - as shown in the screenshot below. We can also see the specific minor number is the source code of the kernel⁴⁹⁴. By the way, the creation of the character device is done as part of the initialization of the module⁴⁹⁵, it's the same “C file” which contains the creation of the “hwrng” kernel thread.

Lastly, to best leverage that we should use “rng-tools”. They use the “/dev/hwrng” to fill the kernel entropy pool,which is used internally and exported by the /dev/urandom and /dev/random special files⁴⁹⁶. By using them we monitor the entropy because just opening and reading from “/dev/hwrng” does not perform any fitness check for the randomness quality⁴⁹⁷.



```
troller$ ls -l /dev/hwrng
crw----- 1 root root 10, 183 07:13 /dev/hwrng
troller$ ls /sys/class/misc/hw_random
dev power rng_available rng_current rng_quality rng_selected subsystem uevent
troller$
```

⁴⁸⁹ https://elixir.bootlin.com/linux/v6.14.5/source/drivers/char/hw_random/core.c#L92
⁴⁹⁰ https://elixir.bootlin.com/linux/v6.14.5/source/drivers/char/hw_random/core.c#L479
⁴⁹¹ <https://medium.com/@boutnaru/the-linux-concept-journey-character-devices-0c75aa70ceb2>
⁴⁹² https://www.kernel.org/doc/Documentation/admin-guide/hw_random.rst
⁴⁹³ <https://medium.com/@boutnaru/the-linux-concept-journey-major-minor-numbers-56abe372482e>
⁴⁹⁴ <https://elixir.bootlin.com/linux/v6.14.5/source/include/linux/miscdevice.h#L46>
⁴⁹⁵ https://elixir.bootlin.com/linux/v6.14.5/source/drivers/char/hw_random/core.c#L670
⁴⁹⁶ https://www.kernel.org/doc/html/next/admin-guide/hw_random.html
⁴⁹⁷ <https://github.com/norman/rng-tools>

psimon (Pressure Stall Information Monitoring)

“psimon” is a kernel thread which is created using the “kthread_create” function⁴⁹⁸. PSI stands for “Pressure Stall Information”, it is used for showing how resource pressure increases as they develop. This is done by leveraging new pressure metrics targeting memory, CPU and IO. Together with other components (cgroup2) we can identify resource shortage while they are developing and react with operations like: pausing\stopping non-essential processes, reallocating memory, load shedding and more⁴⁹⁹.

Overall, PSI is only relevant since kernel version 4.20⁵⁰⁰. The pressure information for each resource (CPU/IO/Memory) is exported to user-mode using profcs: “/proc/pressure/” - as shown in the screenshot below. All of those are created by the “psi_proc_init”⁵⁰¹ function. By the way, the some at the beginning of each line means (in regards to a specific resource) at least one task is delayed on a resource. This affects the workload’s ability to perform work, however the CPU may still be executing other tasks⁵⁰². The full keyword means all non-idle tasks are delayed on that resource such that nobody is advancing and the CPU goes idle. Based on the definition the full state does not exist for the CPU⁵⁰³.

Lastly, for registering a trigger with PSI we just need to open the relevant file related to the resource we want to monitor and write the desired threshold and time window. Then we can leverage select()/poll()/epoll() on the returned file descriptor. An example is “some 170000 2000000 > /proc/pressure/memory” would add 170 ms threshold for partial memory stall measured within 2 sec time window⁵⁰⁴.



```
troller$: cd /proc/pressure/
troller$: ls
cpu io memory
troller$: cat *
some avg10=0.54 avg60=1.69 avg300=1.29 total=430511911
full avg10=0.00 avg60=0.00 avg300=0.00 total=0
some avg10=0.34 avg60=0.61 avg300=0.58 total=127437529
full avg10=0.27 avg60=0.39 avg300=0.40 total=96604458
some avg10=0.00 avg60=0.00 avg300=0.00 total=5064772
full avg10=0.00 avg60=0.00 avg300=0.00 total=4473113
troller$:
```

⁴⁹⁸ <https://elixir.bootlin.com/linux/v6.14.5/source/kernel/sched/psi.c#L1357>

⁴⁹⁹ <https://facebookmicrosites.github.io/psi/docs/overview>

⁵⁰⁰ <https://elixir.bootlin.com/linux/v4.20/source/kernel/sched/psi.c>

⁵⁰¹ <https://elixir.bootlin.com/linux/v4.20/source/kernel/sched/psi.c#L764>

⁵⁰² <https://elixir.bootlin.com/linux/v4.20/source/kernel/sched/psi.c#L26>

⁵⁰³ <https://elixir.bootlin.com/linux/v4.20/source/kernel/sched/psi.c#L30>

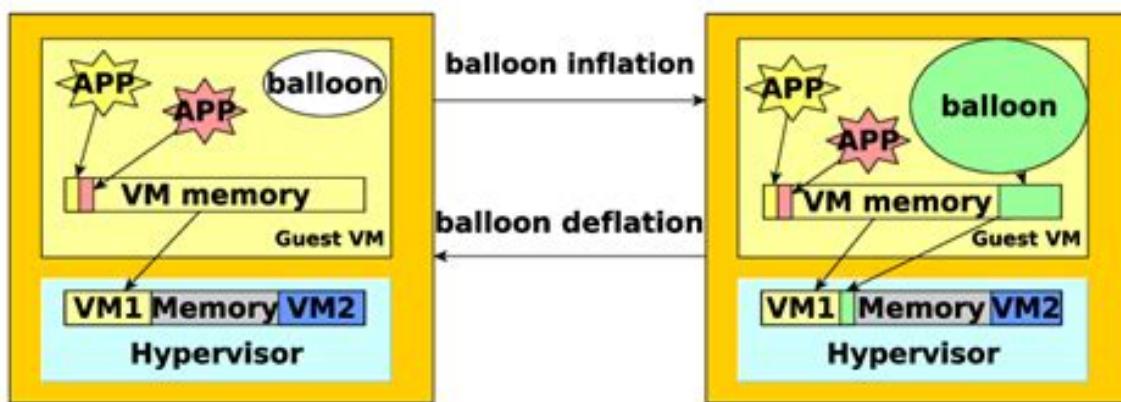
⁵⁰⁴ <https://git.kernel.org/pub/scm/linux/kernel/git/torvalds/linux.git/tree/Documentation/accounting/psi.rst>

hv_balloon (Microsoft Hyper-V Memory Ballooning)

“hv_balloon” is a Linux kernel thread which is created using the “kthread_run” macro⁵⁰⁵, which wraps the “kthread_create” function⁵⁰⁶. The function to run within the thread is “dm_thread_func”⁵⁰⁷ among its operations it posts information regarding the memory pressure to the host every second⁵⁰⁸. “hv_balloon” is part of the Microsoft Hyper-V guest support for Linux⁵⁰⁹, that provides memory ballooning.

Overall, memory ballooning provides the ability of avoiding the need for overcommitting host memory for guest virtual machines (VMs). This is done by letting each VM “give back” unused pages of virtual memory to the host. The balloon driver (inside the guest like “hv_balloon”) allocates unused memory into a pool of memory aka “balloon inflation” (this makes the memory unavailable for applications\processes on the guest). The host (hypervisor\operating system) maps those pages and thus reclaim them aka “balloon deflation”⁵¹⁰ - as shown in the diagram below⁵¹¹.

Lastly, the module name of “Microsoft Hyper-V Balloon driver” is “hv_balloon.ko”. It is part of the Microsoft Hyper-V guest support which is part of the Linux kernel since version 3.1 (release date 2011-10-24). In order to include the support as part of the Linux system we should enable “CONFIG_HYPERV_BALLOON”⁵¹². This support by the way is based on the Hyper-V’s client drivers for Linux⁵¹³.



⁵⁰⁵ https://elixir.bootlin.com/linux/v6.14.6/source/drivers/hv/hv_balloon.c#L1992

⁵⁰⁶ <https://elixir.bootlin.com/linux/v6.14.6/source/include/linux/kthread.h#L42>

⁵⁰⁷ https://elixir.bootlin.com/linux/v6.14.6/source/drivers/hv/hv_balloon.c#L1373

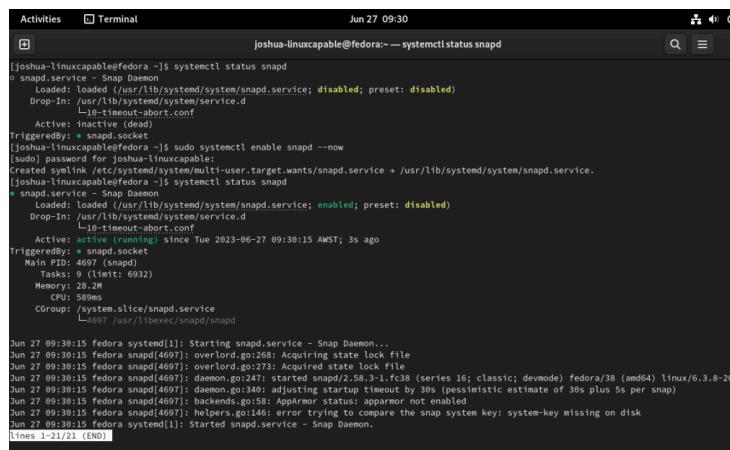
508 https://elixir.bootlin.com/linux/v6.14.6/source/drivers/hv/hv_balloon.c#L1381509 <https://elixir.bootlin.com/linux/v6.14.6/source/drivers/hv/Kconfig#L3>510 https://en.wikipedia.org/wiki/Memory_ballooning511 https://www.researchgate.net/figure/Memory-ballooning-principles_fig2_330462327512 https://www.kernelconfig.io/config_hyperv_balloon513 <https://elixir.bootlin.com/linux/v6.14.6/source/drivers/hv/Kconfig#L51>

snapd (Snap Daemon)

“snapd” is an ELF binary that might be located at “/usr/lib/snapd/snapd”. It is a background service that is used for managing and maintaining snaps⁵¹⁴. Although “snapd” is pre-installed with Ubuntu⁵¹⁵ it can be installed on other Linux distributions⁵¹⁶ such as (but not limited to) Fedora⁵¹⁷ - as shown in the screenshot below.

Overall, “snapd” exposes a REST (Representational State Transfer) API for managing snap⁵¹⁸ packages. We can access those APIs by leveraging the “snap” client which is included in the same package as “snapd”. The snap client can talk with “snapd” by using a Unix domain socket⁵¹⁹ called “snapd.socket”. The socket is commonly located at “/run/snapd.socket”⁵²⁰.

Lastly, “snapd” is an open-source project created by Canonical. It implements confinement policies that isolate snaps from the base system and from each other. Also, it governs the interfaces that allow snaps to access specific system resources outside of their confinement. We can check-out its source code as part of the “snapd” GitHub repository⁵²¹.



```
[joshua-linuxcapable@fedora ~]$ systemctl status snapd
● snapd.service - Snap Daemon
   Loaded: loaded (/usr/lib/systemd/system/snapd.service; disabled; preset: disabled)
     Drop-In: /usr/lib/systemd/system/snapd.service.d
               └─10-timeout-abort.conf
   Active: inactive (dead)
     TriggeredBy: snapd.socket
[joshua-linuxcapable@fedora ~]$ sudo systemctl enable snapd --now
[sudo] password for joshua-linuxcapable:
Created symlink /etc/systemd/system/multi-user.target.wants/snapd.service → /usr/lib/systemd/system/snapd.service.
[joshua-linuxcapable@fedora ~]$ systemctl status snapd
● snapd.service - Snap Daemon
   Loaded: loaded (/usr/lib/systemd/system/snapd.service; enabled; preset: disabled)
     Drop-In: /usr/lib/systemd/system/snapd.service.d
               └─10-timeout-abort.conf
   Active: active (running) since Tue 2023-06-27 09:30:15 AWST; 3s ago
     TriggeredBy: snapd.socket
      Main PID: 469 (snapd)
        Tasks: 1 (limit: 4932)
       Memory: 28.2M
          CPU: 589ms
         CGroup: /system.slice/snapd.service
                   └─469 /usr/libexec/snapd/snapd

Jun 27 09:30:15 fedora systemd[1]: Starting snapd.service - Snap Daemon...
Jun 27 09:30:15 fedora snapd[469]: overlord.go:268: Acquiring state lock file
Jun 27 09:30:15 fedora snapd[469]: overlord.go:273: Acquired state lock file
Jun 27 09:30:15 fedora snapd[469]: daemon.go:247: started snapd/2.58.3-1.fc38 (series 16; classic; devmode) fedora/38 (amd64) linux/6.3.8-20.el8.x86_64
Jun 27 09:30:15 fedora snapd[469]: helpers.go:146: waiting for disk to appear, timeout by 30s (pessimistic estimate of 30s plus 5s per snap)
Jun 27 09:30:15 fedora snapd[469]: backends.go:58: AppArmor status: appears not enabled
Jun 27 09:30:15 fedora snapd[469]: helpers.go:146: error trying to compare the snap system key: system-key missing on disk
Jun 27 09:30:15 fedora systemd[1]: Started snapd.service - Snap Daemon.
Lines 1-21/21 (END)
```

⁵¹⁴ <https://snapcraft.io/docs/installing-snap>

⁵¹⁵ [https://en.wikipedia.org/wiki/Snap_\(software\)](https://en.wikipedia.org/wiki/Snap_(software))

⁵¹⁶ <https://medium.com/@boutnaru/the-linux-concept-journey-linux-distribution-2dfc68aaa4f4>

⁵¹⁷ <https://linuxcapable.com/how-to-install-snap-on-fedora-linux/>

⁵¹⁸ <https://medium.com/@boutnaru/the-linux-concept-journey-snap-f39201df078d>

⁵¹⁹ <https://medium.com/@boutnaru/the-linux-concept-journey-socket-file-0ab8245290d6>

⁵²⁰ <https://wiki.archlinux.org/title/Snap>

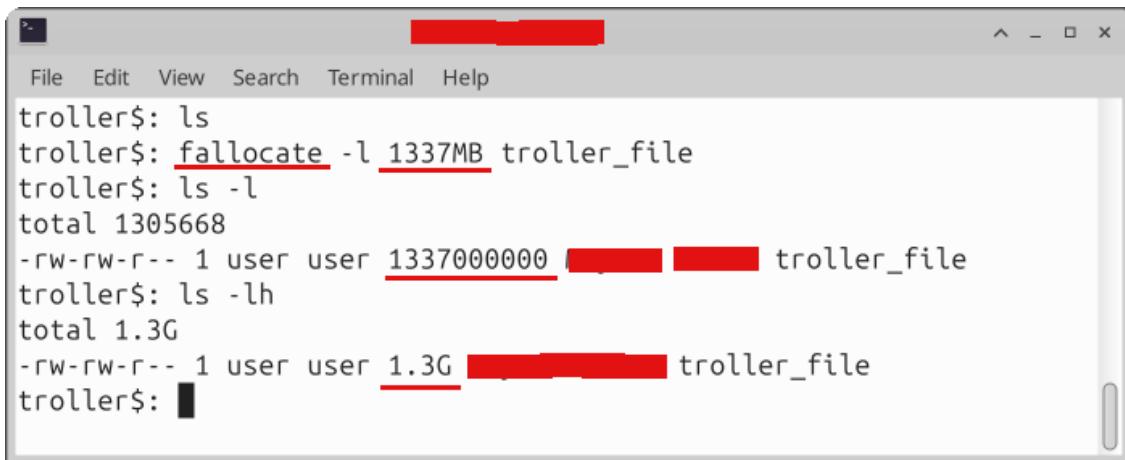
⁵²¹ <https://github.com/canonical/snapd>

fallocate (Preallocate\Deallocate Space for a File)

“fallocate” is an ELF binary located at “/bin/fallocate” and/or “/usr/bin/allocate”. It is used for preallocating and/or deallocating disk space for a file - as shown in the screenshot below. Thus, we can use this command line utility as an alternative for creating a file and filling it with zeros. Due to that users can quickly allocate space without writing any data. It is a great way for creating files with a specific size and/or preventing the case of running out of disk space while writing a file⁵²².

Overall, in case we want to preallocate space on a file-system which natively supports the “fallocate” system call we don’t need to create a file and fill it up with zeros. The preallocation can be done just by allocating blocks and marking them as uninitialized (which might be performed without any IO operations in regards to the data blocks). Thus, it is much more quick than the first option⁵²³.

Lastly, the implementation of the “fallocate” system call as part of libc (using a wrapper) is Linux specific and non-portable. For a POSIX portable use-case we should leverage the “posix_fallocate” library call⁵²⁴. When calling the “fallocate” system call we get into the “vfs_fallocate” that calls the specific file-system code⁵²⁵. We can check the “ntfs_fallocate” function as part of the NTFS code for an implementation example⁵²⁶.



The screenshot shows a terminal window with a gray header bar containing the title 'troller\$'. Below the header is a menu bar with 'File', 'Edit', 'View', 'Search', 'Terminal', and 'Help'. The main area of the terminal contains the following command-line session:

```
troller$: ls
troller$: fallocate -l 1337MB troller_file
troller$: ls -l
total 1305668
-rw-rw-r-- 1 user user 1337000000 [REDACTED] troller_file
troller$: ls -lh
total 1.3G
-rw-rw-r-- 1 user user 1.3G [REDACTED] troller_file
troller$: [REDACTED]
```

⁵²² <https://phoenixnap.com/kb/fallocate>

⁵²³ <https://www.man7.org/linux/man-pages/man1/fallocate.1.html>

⁵²⁴ <https://www.man7.org/linux/man-pages/man2/fallocate.2.html>

⁵²⁵ <https://elixir.bootlin.com/linux/v6.14.6/source/fs/open.c#L338>

⁵²⁶ <https://elixir.bootlin.com/linux/v6.14.6/source/fs/ntfs3/file.c#L561>

ufw (Uncomplicated Firewall)

“ufw” (Uncomplicated Firewall) is a Python script located at “/usr/sbin/ufw” (or “/sbin/ufw”). It is a firewall configuration tool for easing netfilter⁵²⁷ based firewall configuration. Thus, it provides a user-friendly way for creating an IPv4\IPv6 host-based firewall. UFW is the default tool for the job in case of Ubuntu, which is by default disabled. In case we enable ufw (“sudo ufw enable”) it uses a default set of rules (profile) that should be fine for the average home user (or so the developers believe). All ‘incoming’ is denied, with some exceptions to make things easier for home users⁵²⁸.

Overall, UFW is available by default as part of Ubuntu since version 8.04 LTS. Also, it is included as part of Debian since version 10. For easier usage there is also “Gufw” (GUI for Uncomplicated Firewall) which is built using Python, GTK and of course ufw⁵²⁹ - as shown in the screenshot below.

Lastly, “ufw” provides different capabilities such as (but not limited to): default incoming policy (allow/deny), IPv6 support, logging, per rule logging, rsyslog support, rate limiting and filtering by interface. We can checkout the configuration files located at “/etc/ufw/*”⁵³⁰.



⁵²⁷ <https://medium.com/@boutnaru/the-linux-security-journey-netfilter-90c6cf12ca40>

⁵²⁸ <https://help.ubuntu.com/community/UFW>

⁵²⁹ https://en.wikipedia.org/wiki/Uncomplicated_Firewall

⁵³⁰ <https://wiki.ubuntu.com/UncomplicatedFirewall>

arpd (Address Resolution Protocol Daemon)

“arpd” (Address Resolution Protocol Daemon) is an ELF binary located at “/usr/sbin/arpd” or “/sbin/arpd”. It is a user-space arp daemon used for collecting gratuitous ARP messages, which is an ARP response that had been received without sending a request. Gratuitous ARP is a broadcast message announcing\updating the mapping (IP->MAC) of a specific node⁵³¹.

Overall, “arpd” saves the data to disk (by default “/var/lib/arpd/arpd.db ”) and feeds it to the kernel on demand in order to avoid redundant broadcasting - as shown in the screenshot below⁵³². By the way, enabling “CONFIG_ARPD” adds support to have a user space daemon to perform ARP resolution, which is defined in the source code of the Linux kernel up to and including version 3.11.10⁵³³. This configuration basically sets “arpd” as an ARP resolver⁵³⁴.

Lastly, “arpd” supports different switches: “-l” for dumping the database to stdout and “-f” to load\read the database from a file. Also, “-a” sets the number of ARP broadcast queries to send before considering the destination is data; it is often used with “-k” to suppress kernel broadcasts. The “-n” flag configures a timeout for negative cache entries, while “-R” and “-B” control the rate and burst of broadcast packets. We can configure “arpd” to monitor a specific network interface or all of them if none is specified⁵³⁵. For more information we can go over the source code of “arpd” as part of the “iproute2” repository⁵³⁶.

```
root@cloud:~# arpd -b /var/tmp/arpd.db
root@cloud:~# arpd -l -b /var/tmp/arpd.db
#Ifindex IP MAC
2 45 .54 10
root@cloud:~#
```

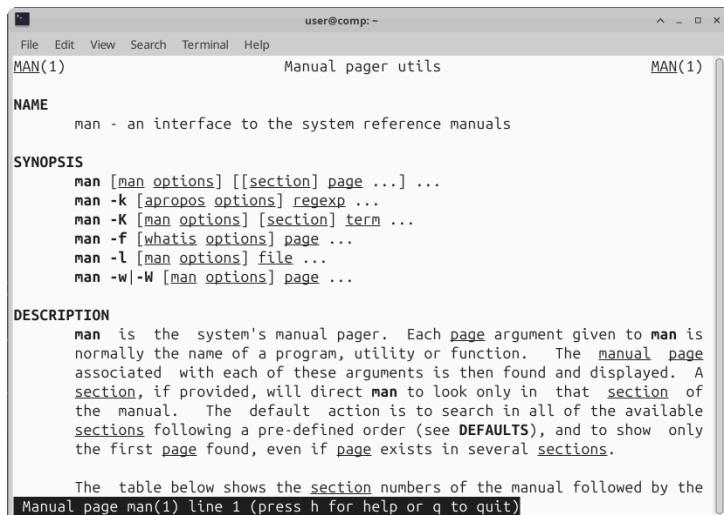
⁵³¹ <https://www.practicalnetworking.net/series/arp/gratuitous-arp/>
⁵³² <https://www.zonghengcloud.com/article/9155.html>
⁵³³ https://elixir.bootlin.com/linux/v3.11.10/A/ident/CONFIG_ARPD
⁵³⁴ <https://linux.die.net/man/8/arpd>
⁵³⁵ <https://man7.org/linux/man-pages/man8/arpd.8.html>
⁵³⁶ <https://github.com/iproute2/iproute2/blob/main/misc/arpd.c>

man (Manual Pages)

man (Manual Pages) is a Linux ELF binary located at “/usr/bin/man” and/or “/bin/man”. It is an interface to the system reference manuals - as shown in the screenshot below. Thus, man is the system’s manual pager⁵³⁷. We can search the manual pages using one of the following equivalent commands: “man -k {expression}”, “man --apropos {expression}” and “apropos {expression}”. By the way, for searching keywords in a whole page text we should use the “-K” option⁵³⁸.

Overall, the Linux man pages are clustered into different topics called “sections”. Among the sections are: “section 1” user commands, “section 2” system calls, “section 3” library calls, “section 4” special files, “section 5” file formats, “section 6” games, “section 7” conventions and miscellany, “section 8” administration and privileged commands, “section L” math library functions and “section N” tcl functions⁵³⁹.

Lastly, “man-db” is an implementation of the standard Unix documentation system, which is leveraged by the “man” utility. It is based on the Berkeley DB database and not the traditional flat-text whatis databases. “man-db” is used by several popular GNU/Linux distributions⁵⁴⁰ such as (but not limited to): Ubuntu, Debian, Arch Linux, Gentoo and OpenSUSE⁵⁴¹. For more information we can check the “man-db” GitLab’s repository⁵⁴².



The screenshot shows a terminal window titled "MAN(1)" with the command "man" entered. The window displays the man(1) manual page for the "man" command. The page is divided into sections: NAME, SYNOPSIS, and DESCRIPTION. The SYNOPSIS section lists various command-line options for "man". The DESCRIPTION section provides a detailed explanation of what "man" does, mentioning it's a manual pager and how it handles multiple sections. A message at the bottom indicates that the table below shows section numbers.

```
user@comp: ~
File Edit View Search Terminal Help
MAN(1)                               Manual pager utils                               MAN(1)

NAME
    man - an interface to the system reference manuals

SYNOPSIS
    man [man options] [[section] page ...] ...
    man -k [apropos options] regexp ...
    man -K [man options] [section] term ...
    man -f [whatis options] page ...
    man -l [man options] file ...
    man -w|-W [man options] page ...

DESCRIPTION
    man is the system's manual pager. Each page argument given to man is
    normally the name of a program, utility or function. The manual page
    associated with each of these arguments is then found and displayed. A
    section, if provided, will direct man to look only in that section of
    the manual. The default action is to search in all of the available
    sections following a pre-defined order (see DEFAULTS), and to show only
    the first page found, even if page exists in several sections.

    The table below shows the section numbers of the manual followed by the
    Manual page man(1) line 1 (press h for help or q to quit)
```

⁵³⁷ <https://man7.org/linux/man-pages/man1/man.1.html>

⁵³⁸ https://wiki.archlinux.org/title/Man_page

⁵³⁹ <https://linux.die.net/man/>

⁵⁴⁰ <https://medium.com/@boutinraru/the-linux-concept-journey-linux-distribution-2dfc68aaa4f4>

⁵⁴¹ <https://man-db.nongnu.org/>

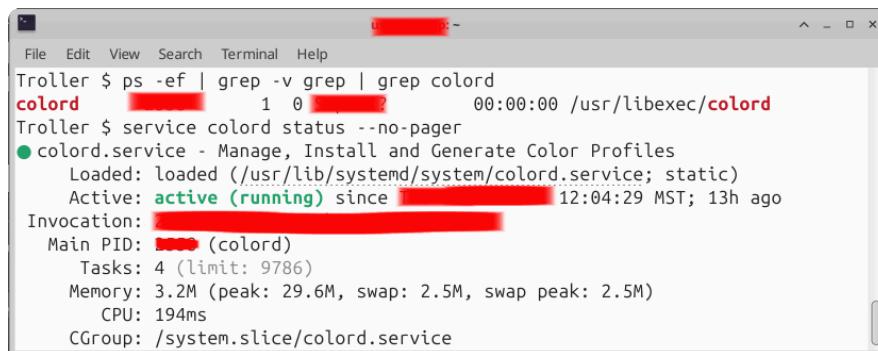
⁵⁴² <https://gitlab.com/man-db/man-db>

colord (Color Daemon)

“colord” is an ELF binary that can be located at “/usr/libexec/colord”. It is a system service which provides the ability to install\generate color profiles (this is for color-managed input and output devices). It is used by “GNOME Color Manager” for system integration. It is also used when no user is logged in⁵⁴³. By default, the “colord” service is executed with the permissions of the “colord” user - as shown in the screenshot below.

Overall, “colord” provides different capabilities such as (but not limited to): persistent database backed store that is preserved across reboots, a D-Bus⁵⁴⁴ API for system frameworks to query and provides the ability to set color based system settings. “colord” supports various subsystems like: XRandR (monitors), SANE (scanners), UDEV (cameras), CUPS (printers) and Virtaul (scanners/cameras/printers). Both “GNOME Color Manager” and “colord-kde” (used with KDE desktop) act as clients of “colord”⁵⁴⁵.

Lastly, colord leverages PolKit⁵⁴⁶ for user authentication. Thus, we can configure what users can do and what not⁵⁴⁷. The configuration of devices to profiles is stored in the colord’s mapping database (/var/lib/colord/mapping.db). Persistent devices\profiles can also be saved to “/var/lib/colord/storage.db” aka the storage database⁵⁴⁸. For more information I suggest going over colord’s GitHub repository⁵⁴⁹.



The screenshot shows a terminal window with the following content:

```
Troller $ ps -ef | grep -v grep | grep colord
colord 1 0 0:00:00 /usr/libexec/colord
Troller $ service colord status --no-pager
● colord.service - Manage, Install and Generate Color Profiles
  Loaded: loaded (/usr/lib/systemd/system/colord.service; static)
  Active: active (running) since [REDACTED] 12:04:29 MST; 13h ago
    Invocation: [REDACTED]
      Main PID: [REDACTED] (colord)
        Tasks: 4 (limit: 9786)
       Memory: 3.2M (peak: 29.6M, swap: 2.5M, swap peak: 2.5M)
         CPU: 194ms
        CGroup: /system.slice/colord.service
```

⁵⁴³ <https://www.linuxfromscratch.org/blfs/view/svn/general/colord.html>

⁵⁴⁴ <https://medium.com/@boutnaru/the-linux-concept-journey-d-bus-desktop-bus-dd8c69ade019>

⁵⁴⁵ <https://www.freedesktop.org/software/colord/intro.html>

⁵⁴⁶ <https://medium.com/@boutnaru/the-linux-security-journey-polkit-authorization-manager-459c2dd2c3cf>

⁵⁴⁷ <https://www.freedesktop.org/software/colord/species.html>

⁵⁴⁸ <https://www.freedesktop.org/software/colord/species.html>

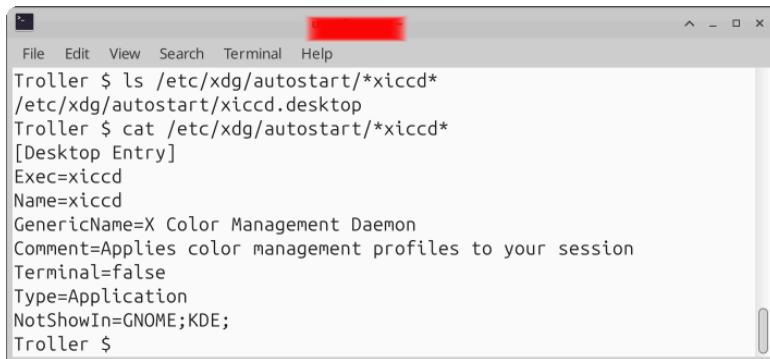
⁵⁴⁹ <https://github.com/huehsie/colord>

xiccd (X Color Management Daemon)

xiccd (X Color Management Daemon) is an ELF binary located at “/usr/bin/xiccd” and/or “/bin/xiccd”. It is basically a bridge between colord⁵⁵⁰ and X⁵⁵¹. It can be used for performing various operations such as (but not limited to): enumerating displays, registering displays in “colord”, creating default ICC (International Color Consortium) profiles based on EDID (Extended Display Identification Data), applying ICC profiles and maintaining private ICC storage directory⁵⁵².

Overall, it allows non-GNOME and non-KDE desktop environments to load and apply icc profiles⁵⁵³. By the way, an ICC profile is set of data that characterizes a color input or output device, or a color space⁵⁵⁴. xiccd is not dependent on any specific desktop environment\toolkit. Also, it should not be used in desktop environments that support color management natively, for example GNOME and KDE⁵⁵⁵.

Lastly, one example of a desktop environment is xfce that does not support applying display profiles without xiccd⁵⁵⁶. xiccd is typically launched automatically on login like by leveraging “/etc/xdg/autostart”⁵⁵⁷ - as shown in the screenshot below. For more information I suggest going over xiccd’s code in its GitHub repository⁵⁵⁸.



The screenshot shows a terminal window titled 'Troller' with a redacted title bar. The window has a standard Linux-style menu bar with File, Edit, View, Search, Terminal, and Help. The terminal content shows the user navigating to /etc/xdg/autostart and listing files related to xiccd. Then, the user cat'ing the xiccd.desktop file. The file's contents are displayed:

```
File Edit View Search Terminal Help
Troller $ ls /etc/xdg/autostart/*xiccd*
Troller $ cat /etc/xdg/autostart/*xiccd*
[Desktop Entry]
Exec=xiccd
Name=xiccd
GenericName=X Color Management Daemon
Comment=Applies color management profiles to your session
Terminal=false
Type=Application
NotShowIn=GNOME;KDE;
Troller $
```

⁵⁵⁰ <https://medium.com/@boutnaru/the-linux-process-journey-colord-color-daemon-9e718013c679>

⁵⁵¹ <https://medium.com/@boutnaru/the-linux-concept-journey-xorg-0d7ab9697a34>

⁵⁵² <https://manpages.debian.org/trixie/xiccd/xiccd.8.en.html>

⁵⁵³ https://wiki.archlinux.org/title/ICC_profiles

⁵⁵⁴ https://en.wikipedia.org/wiki/ICC_profile

⁵⁵⁵ <https://community.linuxmint.com/software/view/xiccd>

⁵⁵⁶ <https://docs.xfce.org/xfce/xfce4-settings/color>

⁵⁵⁷ <https://discussion.fedoraproject.org/t/abhis3k-xiccd/83440>

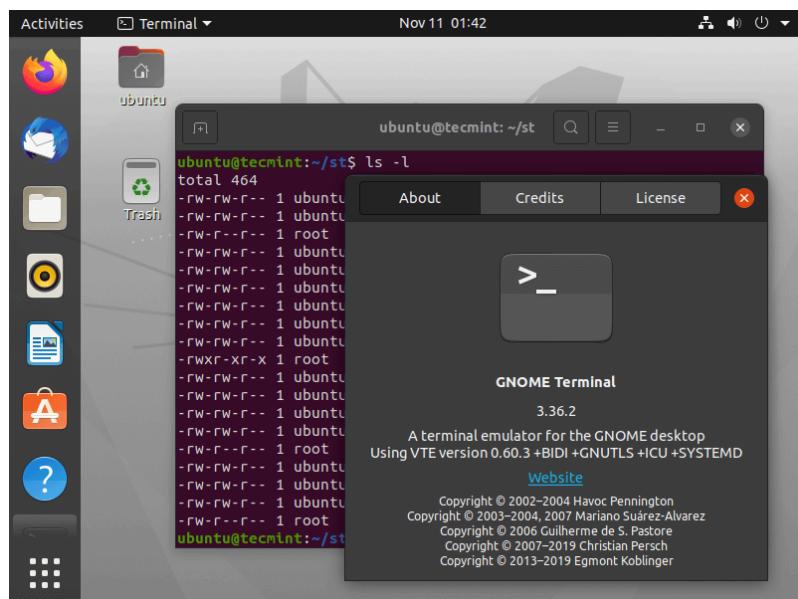
⁵⁵⁸ <https://github.com/aalakhov/xiccd>

gnome-terminal

“gnome-terminal” is an ELF binary located at “/usr/bin/gnome-terminal” and/or “//bin/gnome-terminal”. It is a terminal emulator application which can be leveraged for accessing Linux shell environments. For example running programs available on your system⁵⁵⁹ - as shown in the screenshot below⁵⁶⁰.

Overall, it supports different capabilities such as (but not limited to): using tabs, executing commands, working with text (copy/paste/search/selection/etc), customizing appearance (changing cursor/setting terminal size/zooming in/zooming out/color selection/etc), customizing behaviour (disabling user input/manage profiles/keyboard accessibility/etc) and resetting terminal state⁵⁶¹.

Lastly, “gnome-terminal” has multiple command line switches like which can be clustered to different groups: general information (like “-h”\”--help” for overview of all options), window management (like “--tab” which opens a new tab), executing commands (“--command”\”-c”), appearance management (like “-t”\”--title” for setting the terminal title) and others⁵⁶². For more information I suggest going over the code repository of “gnome-terminal”⁵⁶³.



⁵⁵⁹ <https://help.gnome.org/users/gnome-terminal/stable/introduction.html.en>

⁵⁶⁰ <https://www.tecmint.com/linux-terminal-emulators/>

⁵⁶¹ <https://help.gnome.org/users/gnome-terminal/stable/>

⁵⁶² <https://gemini.google.com/app/c557a41018e27760>

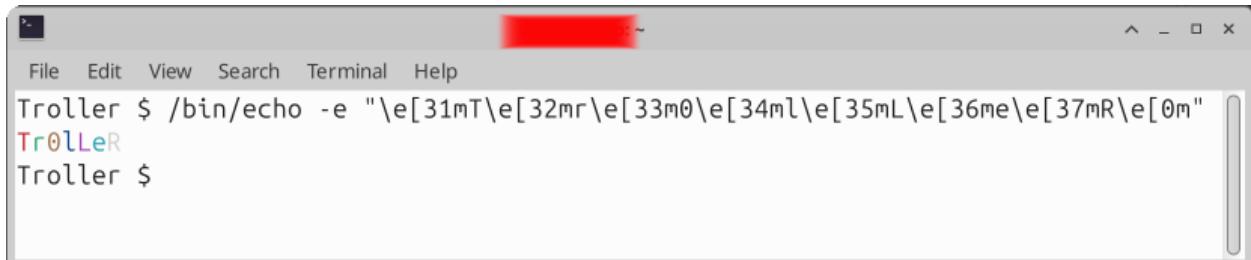
⁵⁶³ <https://gitlab.gnome.org/GNOME/gnome-terminal>

echo

“echo” is an ELF binary located at “/usr/bin/echo” and/or “/bin/echo”. Its simplest usage form is to display a line of text⁵⁶⁴. echo also supports specific escaped characters for different use-cases like (but not limited to): horizontal tab (“\t”), vertical tab (“\v”), “new line” (“\n”), remove spaces (“\b”) and the usage of backslash itself (“\\”). Also, when using bash (and other shells) we can leverage color displaying using the following pattern: “\e[<STYLE>;<COLOR>m”⁵⁶⁵ - as shown in the screenshot below.

Overall, we can use different colors such as (but not limited to): black, green, yellow, blue and purple. Moreover, there is also support for underline, bold, background color and more⁵⁶⁶. By the way, “echo” is part of the coreutils package⁵⁶⁷. Hence, we can also check the source code of “echo” as part of the coreutils’ GitHub repository⁵⁶⁸.

Lastly, it is important to understand that shells like bash can have a builtin implementation of “echo”. In that case invoking it without the full path causes the internal implementation to run⁵⁶⁹. There are also equivalent implementations for other operating systems like Windows⁵⁷⁰.



The screenshot shows a terminal window with a redacted title bar. The menu bar includes File, Edit, View, Search, Terminal, and Help. The command entered is `/bin/echo -e "\e[31mT\e[32mr\e[33m0\e[34ml\e[35mL\e[36me\e[37mR\e[0m"`. The output displayed is `Tr0lLeR`, where each character has a different color: T is red, r is green, 0 is yellow, l is blue, L is magenta, e is cyan, and R is black. The prompt `Troller $` is visible at the bottom.

⁵⁶⁴ <https://linux.die.net/man/1/echo>

⁵⁶⁵ <https://ss64.com/bash/echo.html>

⁵⁶⁶ <https://gist.github.com/lomedil/41b739a74b481c4b0a47fca09f42bea3>

⁵⁶⁷ <https://www.gnu.org/software/coreutils/>

⁵⁶⁸ <https://github.com/coreutils/coreutils/blob/master/src/echo.c>

⁵⁶⁹ https://www.gnu.org/software/bash/manual/html_node/Bash-Builtins.html

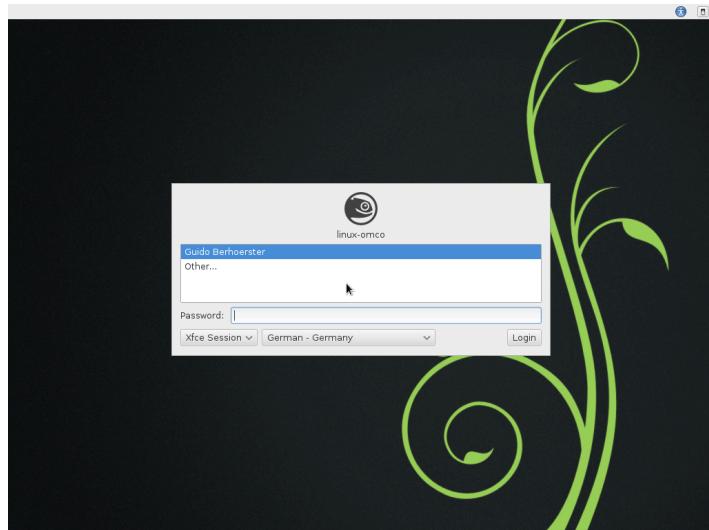
⁵⁷⁰ <https://ss64.com/nt/echo.html>

lightdm

LightDM is an ELF binary located at “/usr/sbin/lightdm” and\or “/sbin/lightdm”. It is used as a cross-desktop display manager for Linux based systems⁵⁷¹. LightDM was developed (at Canonical) as a relatively light-weight and highly customizable alternative to GDM (GNOME Display Manager)⁵⁷². The configuration file is usually stored at “/etc/lightdm/lightdm.conf”. By the way, we can create “/etc/lightdm/lightdm.conf.d/” and place our configuration files there⁵⁷³.

Overall, among LightDM features are (but not limited to): supporting various desktop environments, supporting different display technologies (like X\Wayland\Mir), low memory usage, supporting guest sessions, supporting remote logins (incoming XDMCP\VNC and outgoing XDMCP\PAM) and a comprehensive test suite. Moreover, if we want a GUI that prompts the user for credentials\lets the user select a session\etc we need to install a greeter - as shown below⁵⁷⁴ leveraging gtk greeter. It is possible to use LightDM without a greeter only in case of an automatic login⁵⁷⁵.

Lastly, by default it is executed with the permissions of the root user⁵⁷⁶. For more information I suggest going over LightDM's GitHub repository. Based on the repo's statistics it is ~89% of the source code is written in the C programming language⁵⁷⁷. As opposed to GDM it does not load any GNOME libraries⁵⁷⁸.



⁵⁷¹ <https://medium.com/@boutnaru/the-linux-concept-journey-display-manager-dm-05634c3e8c29>

⁵⁷² <https://wiki.debian.org/LightDM>

⁵⁷³ <https://wiki.debian.org/LightDM>

574 <https://en.opensuse.org/File:Lightdm-display-manager-gtk-greeter.png>
575 <https://en.opensuse.org/File:Lightdm-gtk-greeter-11.1-LightDM.png>

<https://wiki.archlinux.org/title/LightDM>

<https://man.archlinux.org/man/lightdm.1.en>

<https://en.wikipedia.org/wiki/LightDM>

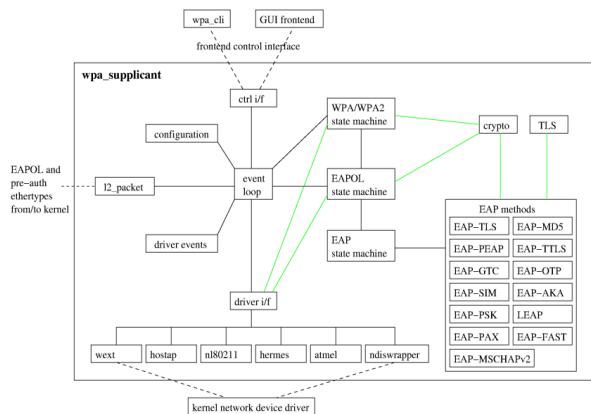
<https://en.wikipedia.org/wiki/LightDM>

wpa_supplicant (Wi-Fi Protected Access client)

“wpa_supplicant” is an ELF binary located at “/usr/sbin/wpa_supplicant” and/or “/sbin/wpa_supplicant”. It is a cross-platform supplicant (a client device that needs to be verified by a server before using a dedicated network link) with support for WPA\WPA2\WPA3 (IEEE 802.11i). Thus, it is used by desktops\laptops\mobile devices\embedded systems. It implements key negotiation with a WPA authenticator and it controls the roaming and IEEE 802.11 authentication/association of the wireless driver⁵⁷⁹.

Overall, “wpa_supplicant” is meant for executing as a “daemon” (user-mode) in the background and acting as the backend component controlling the wireless connection. Hence, it will run with the permissions of the “root” user. It is important to know that “wpa_supplicant” has also frontend programs “wpa_cli” (text-based frontend) and “wpa_gui” (GUI based)⁵⁸⁰. As an example the phases associating with an AP using WPA are as follows. First, requests the kernel driver to scan neighboring BSSes and then selects and associates with one. If using WPA-EAP, it completes EAP authentication to get the master key; if using WPA-PSK, it uses the PSK as the master key. Next, completing the WPA 4-Way and Group Key Handshakes with the AP. In the last phase, it configures the encryption keys, allowing for normal data transmission⁵⁸¹.

Lastly, “wpa_supplicant” supports different configurations stored by default at “/etc/wpa_supplicant/”. Probably the most fundamental one is “wpa_supplicant.conf”⁵⁸². For more details about the configuration file I suggest going over a detailed one as an example⁵⁸³. Also, “wpa_supplicant” exposes a control interface for providing the ability of external programs to control the operations of the daemon and to get status information and event notifications. For more information I suggest going over the source code⁵⁸⁴. A general map of the modules related to “wpa_supplicant” is shown below⁵⁸⁵.



⁵⁷⁹ https://wiki.archlinux.org/title/Wpa_supplicant

⁵⁸⁰ https://w1.fi/wpa_supplicant/

581 https://linux.die.net/man/8/wpa_supplicant

https://wiki.gentoo.org/wiki/Wpa_supplicant

https://git.w1.fi/cgit/hostap/tree/wpa_supplicant/wpa_supplicant.conf

⁵⁸⁴ <https://git.w1.fi/cgit/hostap/>

585 https://w1.fi/wpa_supplicant-devel/

wpa_cli (Wi-Fi Protected Access Command Line Client)

“wpa_cli” (Wi-Fi Protected Access Command Line Client) is an ELF binary located at “/usr/sbin/wpa_cli” and/or “/sbin/wpa_cli”. It is a text-based frontend program for which interacts with the user-mode daemon “wpa_supplicant”⁵⁸⁶. We can leverage it for querying current status\changing configuration\triggering events and requesting interactive user input⁵⁸⁷.

Overall, “wpa_cli” provides an interactive mode which supports tab compilation and description of complemented commands - as shown in the screenshot below. By default, “wpa_cli” uses the control socket located at “/var/run/wpa_supplicant/”. It is important to know that for running “wpa_cli” with all of its features we need to have root permissions or the relevant capabilities⁵⁸⁸.

Lastly, for more information about the commands supported by “wpa_cli” and their arguments I suggest going over different cheatsheet⁵⁸⁹. Moreover, “wpa_cli” is packaged as part of the “wpa_supplicant” in most Linux based operating systems⁵⁹⁰. In order to get a better understanding about “wpa_cli” we can go over its source code⁵⁹¹.

```
> scan
OK
<3>CTRL-EVENT-SCAN-RESULTS
> scan_results
bssid / frequency / signal level / flags / ssid
00:00:00:00:00:00 2462 -49 [WPA2-PSK-CCMP][ESS] MYSSID
11:11:11:11:11:11 2437 -64 [WPA2-PSK-CCMP][ESS] ANOTHERSSID
```

⁵⁸⁶ <https://medium.com/@boutnaru/the-linux-process-journey-wpa-supplicant-wi-fi-protected-access-client-eef744a71e00>

⁵⁸⁷ https://linux.die.net/man/8/wpa_cli

⁵⁸⁸ https://wiki.archlinux.org/title/Wpa_supplicant

⁵⁸⁹ <https://gist.github.com/penguinpownz/1d36a38af4fac4553562410e0bd8d6cf>

⁵⁹⁰ <https://packages.debian.org/trixie/amd64/wpasupplicant/filelist>

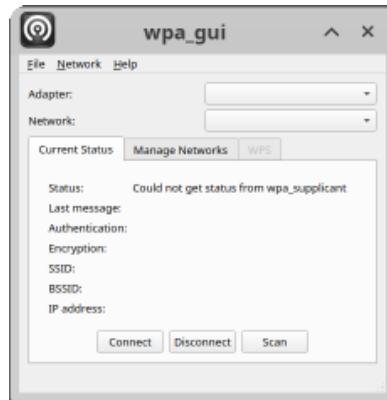
⁵⁹¹ https://android.googlesource.com/platform/external/wpa_supplicant/+/6eb6225dda14c6eaf34fcb3498e17d28060e1e16/wpa_cli.c

wpa_gui (Wi-Fi Protected Access Graphical User Interface)

“wpa_gui” (Wi-Fi Protected Access Graphical User Interface) is an ELF binary. We can find it at “/usr/sbin/wpa_gui” and/or “/sbin/wpa_gui”. However, those locations can be a symbolic link to “/usr/bin/wpa_gui” and/or “/bin/wpa_gui”. It is a GUI based frontend program for which interacts with the user-mode daemon “wpa_supplicant”⁵⁹² - as shown below. Thus, we can use “wpa_gui” instead of NetworkManager⁵⁹³ for configuring and managing wireless interfaces⁵⁹⁴.

Overall, “wpa_gui” is based on QT which supports almost all the interactive statuses\configuration features of “wpa_cli”⁵⁹⁵. It also uses that default location for “wpa_supplicant” that is “/var/run/wpa_supplicant/”. It can be overridden using the “-p” command line switch of “wpa_gui”⁵⁹⁶.

Lastly, on some Linux distributions (like Ubuntu) “wpa_gui” is included in a different package from “wpa_cli” and “wpa_supplicant” called “wpagui”⁵⁹⁷. For more information I suggest going over an implementation of “wpagui”⁵⁹⁸.



⁵⁹² <https://medium.com/@boutnaru/the-linux-process-journey-wpa-suplicant-wi-fi-protected-access-client-cef744a71e00>

⁵⁹³ <https://medium.com/@boutnaru/the-linux-process-journey-modemmanager-modem-management-daemon-71d3b5cf9db1>

⁵⁹⁴ <https://ansible-linux-postinstall.readthedocs.io/en/latest/guide-task-wpagui.html>

⁵⁹⁵ <https://medium.com/@boutnaru/the-linux-process-journey-wpa-cli-wi-fi-protected-access-command-line-client-931dc8b28ffb>

⁵⁹⁶ https://linux.die.net/man/8/wpa_gui

⁵⁹⁷ <https://launchpad.net/ubuntu/noble/+package/wpagui>

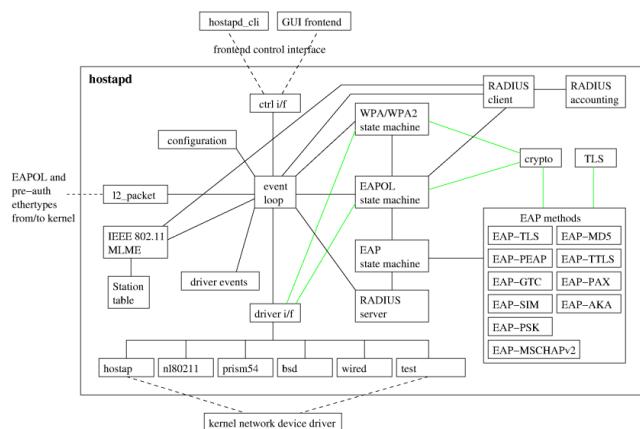
⁵⁹⁸ https://android.googlesource.com/platform/external/wpa_supplicant/+refs/heads/tools_r20/wpa_gui/

hostapd (Host Access Point Daemon)

“hostapd” (Host Access Point Daemon) is an ELF binary located at “/usr/sbin/hostapd” and/or “/sbin/hostapd”. It is a user-mode daemon used as an AP (access point) and an authentication server. Among the protocols which are supported are: IEEE 802.1X/WPA/EAP/RADIUS⁵⁹⁹. By the way, in order to communicate with a kernel driver it leverages different interfaces. All new cfg80211 (and mac80211) based drivers that implement AP functionality are supported using the nl80211 interface. For old kernel drivers it has drivers, there are 3 other drivers you can use: “madwifi”, “prism54” and “HostAP”⁶⁰⁰.

Overall, hostapd provides various capabilities such as (but not limited to): creating an AP, creating multiple APs on the same card (if the card supports it), create different APs on different interfaces using the same instance of the daemon and using both 2.4GHz and 5GHz at the same time on the same card (requires a support from the hardware side also). It is important to understand that there are features not implemented\supported by hostapd like: assigning an IP address to the AP itself, assigning IPs to the devices connecting to the AP (DHCP server is needed), creating a dual-band AP, even with two cards (it can create two APs with the same SSID) and creating multiple APs on different channels (on the same card). Multiple APs on the same card will share the same channel⁶⁰¹.

Lastly, the source code of “hostapd” is divided into separate C files - as shown below. All hardware/driver specific functionality is in separate files that implement a well-defined driver AP⁶⁰². There is also specific information regarding porting to different target boards and operating systems⁶⁰³. For more information I suggest going over the source code of hostapd⁶⁰⁴.



⁵⁹⁹ <https://w1.fi/>

⁶⁰⁰ <https://wireless.docs.kernel.org/en/latest/en/users/documentation/hostapd.html>

⁶⁰¹ <https://wiki.gentoo.org/wiki/Hostapd>

⁶⁰² https://w1.fi/wpa_supplicant-devel/

⁶⁰³ https://w1.fi/wpa_supplicant-devel/porting.html

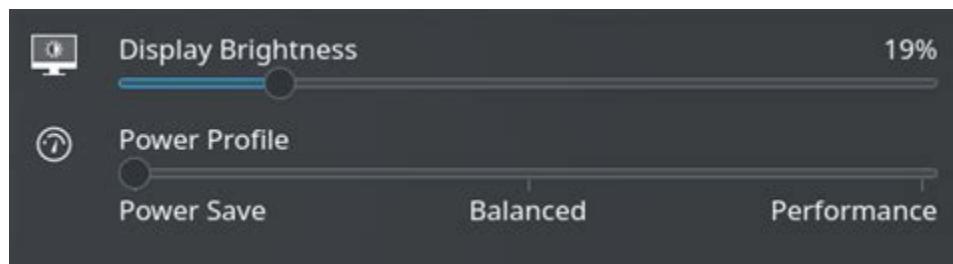
⁶⁰⁴ https://w1.fi/wpa_supplicant-devel/dir_68267d1309a1af8e8297ef4c3efbcd.html

power-profiles-daemon

“power-profiles-daemon” is an ELF binary located by default at “/usr/libexec/power-profiles-daemon”. It is used for modifying the system power/behavior state. Thus, it is leveraged on various laptops\desktop environments to activate power saving and\or performance CPU governors through D-BUS⁶⁰⁵. “power-profiles-daemon” is used as a tool for setting the CPU governor in order to increase system performance at the cost of energy usage⁶⁰⁶. It is mostly executed with the permissions of the “root” user.

Overall, there are three different power profiles that a user can select from: "balanced" (default mode), "power-saver" mode and "performance" mode. The first two are available on every system while "performance" is available on selected systems only (implemented by different "drivers"). By the way, we can also hook to specific changes in the behaviour of a particular device. For example, we can disable fast-charging for some USB devices in power-saver mode⁶⁰⁷.

Lastly, the selection of profiles can be done by using different settings as part of the desktop environment - as shown in the screenshot below from Arch Linux running KDE⁶⁰⁸. Hence, we can summarize that “power-profiles-daemon” is a CPU throttle that is most effective when high\medium load applications are in use⁶⁰⁹. For more information I suggest going over the code repository of “power-profiles-daemon”⁶¹⁰.



⁶⁰⁵ <https://medium.com/@boutnaru/the-linux-concept-journey-d-bus-desktop-bus-dd8c69ade019>

⁶⁰⁶ <https://www.linuxfromscratch.org/blfs/view/systemd/general/power-profiles-daemon.html>

⁶⁰⁷ <https://github.com/Rongronggg9/power-profiles-daemon>

⁶⁰⁸ https://www.reddit.com/r/kde/comments/r74fp/change_cpu_governor_using_powerprofilesdaemon/

⁶⁰⁹ <https://linrunner.de/tlp/faq/ppd.html>

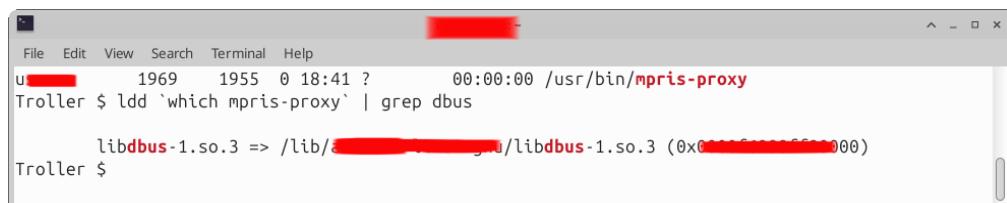
⁶¹⁰ <https://gitlab.freedesktop.org/u-power/power-profiles-daemon>

mpris-proxy (Media Player Remote Interfacing Specification Proxy)

“mpris-proxy” (Media Player Remote Interfacing Specification Proxy) is an ELF binary located at “/usr/bin/mpris-proxy” and\or “/bin/mpris-proxy”. It is part of the “bluez” package⁶¹¹. MPRIS is a standard D-Bus⁶¹² interface which is used for providing a common API for controlling media players - as shown below. This API is leveraged for different tasks such as (but not limited to) discovery and playback⁶¹³.

Overall, by using “mpris-proxy” we enable a proxy forwarding Bluetooth MIDI (Musical Instrument Digital Interface) controls via MPRIS to control media players⁶¹⁴. The process of the “mpris-proxy” is executed with the permissions of the logged on user - as shown below.

Lastly, we can summarize the functionality of “mpris-proxy” to the following: registering media players, interface translation (BlueZ<->MPRIS), volume management, cover art support and more . For more information I suggest going over the source code of “mpris-proxy” as part of the “bluez” GitHub’s repository⁶¹⁵.



```
File Edit View Search Terminal Help
u 1969 1955 0 18:41 ? 00:00:00 /usr/bin/mpris-proxy
Troller $ ldd `which mpris-proxy` | grep dbus
libdbus-1.so.3 => /lib/canonical/lib/libdbus-1.so.3 (0x0000000000000000)
Troller $
```

⁶¹¹ <https://askubuntu.com/questions/1349701/bluetooth-speaker-or-headphones-volume-and-playback-control-mpris>

⁶¹² <https://medium.com/@boutnaru/the-linux-concept-journey-d-bus-desktop-bus-dd8c69ade019>

⁶¹³ <https://wiki.archlinux.org/title/MPRIS>

⁶¹⁴ <https://mynixos.com/home-manager/options/services.mpris-proxy>

⁶¹⁵ <https://github.com/bluez/bluez/blob/master/tools/mpris-proxy.c>