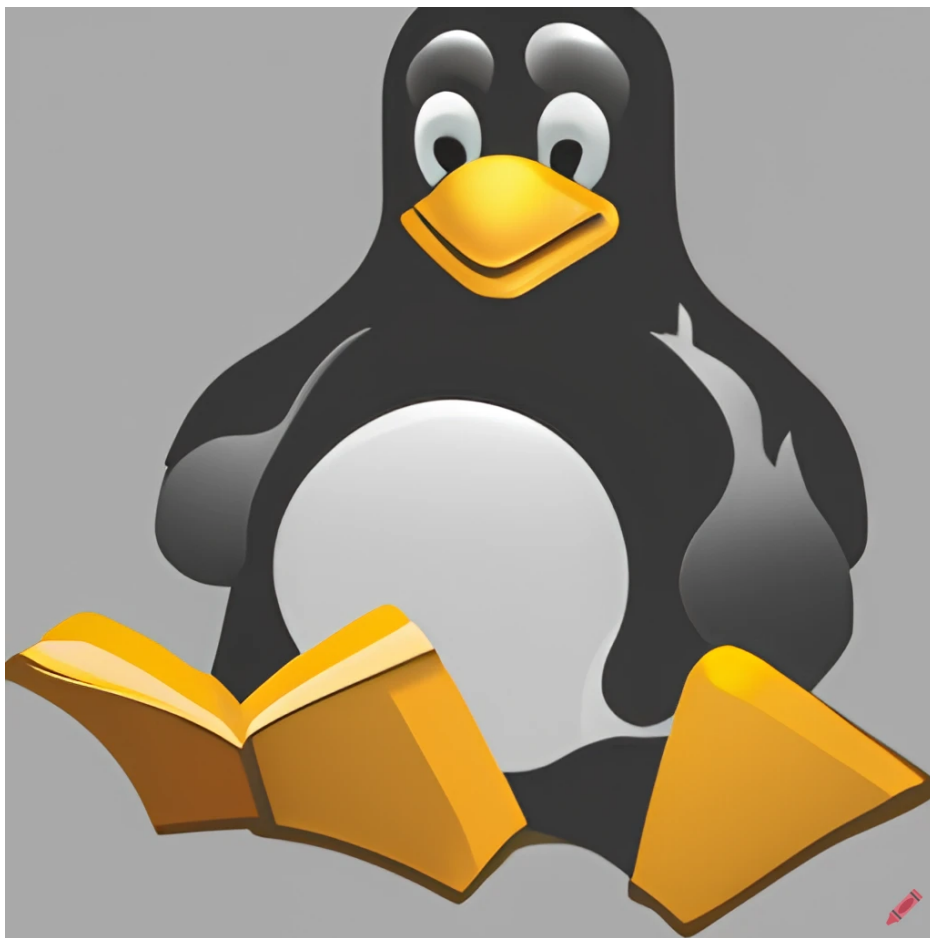


# The Linux Process Journey

version 3.0 (beta)

May-2023

By Dr. Shlomi Boutnaru



Created using [Craiyon AI Image Generator](#)

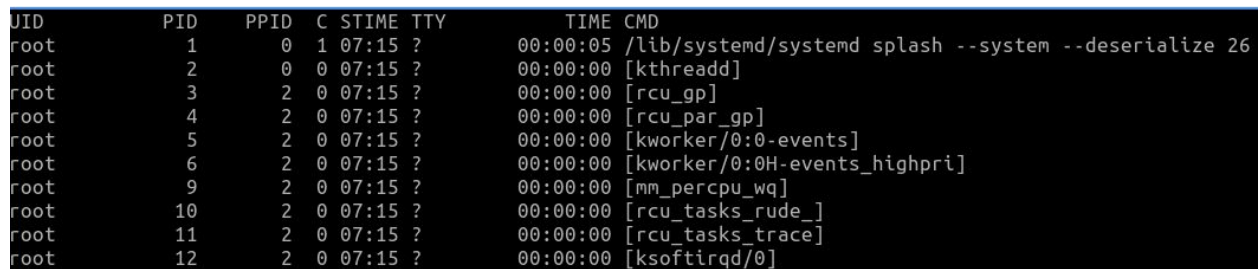
<b>Introduction</b>	<b>3</b>
<b>swapper (PID 0)</b>	<b>4</b>
<b>init (PID 1)</b>	<b>5</b>
<b>Kernel Threads</b>	<b>6</b>
<b>kthreadd (PID 2)</b>	<b>7</b>
<b>migration</b>	<b>9</b>
<b>charger-manager</b>	<b>11</b>
<b>idle_inject</b>	<b>12</b>
<b>kworker</b>	<b>13</b>
<b>kdevtmpfs</b>	<b>14</b>
<b>cpuhp</b>	<b>16</b>
<b>khungtaskd</b>	<b>16</b>
<b>kswapd</b>	<b>18</b>
<b>kcompactd</b>	<b>19</b>
<b>md (Multiple Device Driver)</b>	<b>21</b>
<b>mld (Multicast Listener Discovery)</b>	<b>23</b>
<b>ksmd (Kernel Same Page Merging)</b>	<b>24</b>
<b>ttm_swap</b>	<b>26</b>
<b>watchdogd</b>	<b>28</b>
<b>zswap-shrink</b>	<b>30</b>
<b>khugepaged</b>	<b>31</b>
<b>krfcommd</b>	<b>32</b>
<b>ksgxd</b>	<b>33</b>
<b>jbd2 (Journal Block Device 2)</b>	<b>34</b>
<b>netns</b>	<b>35</b>
<b>oom_reaper</b>	<b>36</b>
<b>kpsmoused</b>	<b>37</b>
<b>slub_flushwq</b>	<b>38</b>
<b>pgdatinit</b>	<b>39</b>
<b>kblockd</b>	<b>40</b>
<b>writeback</b>	<b>41</b>
<b>kdamond (Data Access MONitor)</b>	<b>42</b>
<b>kintegrityd</b>	<b>43</b>
<b>kthrotld</b>	<b>44</b>
<b>scsi_eh (Small Computer System Interface Error Handling)</b>	<b>45</b>
<b>blkcg_punt_bio</b>	<b>46</b>

# Introduction

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

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

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



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

Probably the best way to do it is to go over the processes by the order of their PID value.

The first one I want to talk about is the one we can't see on the list, that is PID 0 (we can see it is the PPID for PID 1 and PID 2 - on them in the next posts).


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

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



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

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

# init (PID 1)

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

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

The creation of init is done by the kernel function “rest\_init”<sup>1</sup>. In the code we can see the call to “user\_mode\_thread” which spawns init, later in the function there is a call to “kernel\_thread” which creates PID 2 (more information about it in the upcoming pages ;-).

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

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

---

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

# Kernel Threads

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

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

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

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

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

## kthreadd (PID 2)

After explaining about PID 1, now we are going to talk about PID 2. Basically, kthreadd is the “kernel thread daemon”. Creation of a new kernel thread is done using kthreadd (We will go over the entire flow). Thus, the PPID of all kernel threads is 2 (checkout ps to verify this). As explained in the post about PID 1 (init) the creation of “kthreadd” is done by the kernel function “rest\_init”<sup>2</sup>. 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 loop<sup>3</sup>) and checks “kthread\_create\_list” for new kthreads to be created. In order to create a kthread the function “kthread\_create”<sup>4</sup> is used, which is a helper macro for “kthread\_create\_on\_node”<sup>5</sup>. We can also call “kthread\_run”<sup>6</sup> 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”<sup>7</sup>). 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.

---

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

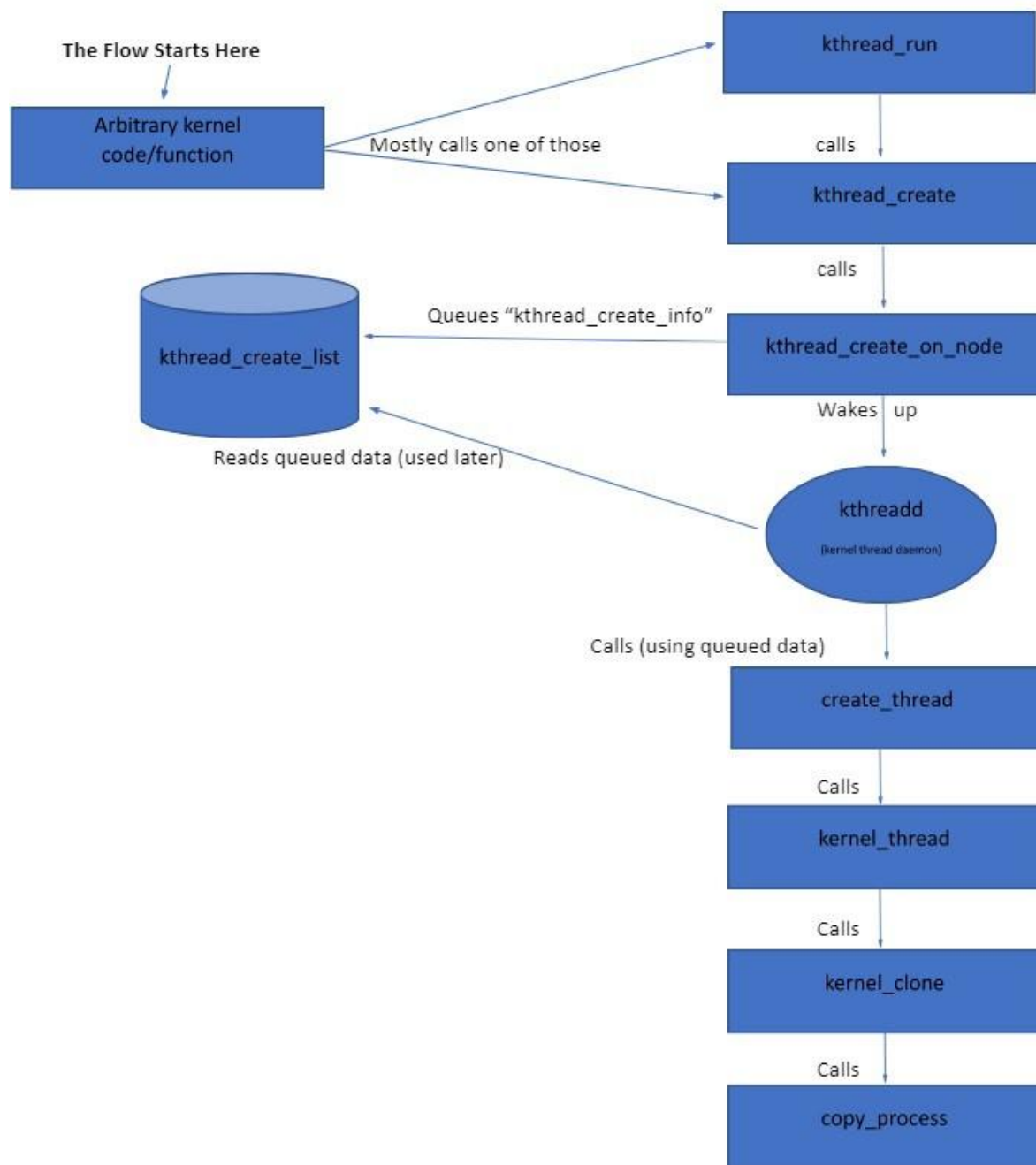
<sup>3</sup> <https://elixir.bootlin.com/linux/v6.1.12/source/kernel/kthread.c#L731>

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

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

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

<sup>7</sup> <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()”<sup>8</sup>). 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-recover to 0 CPU
migration/2 moved lvm2-activation to 0 CPU
migration/0 moved (direxec) to 2 CPU
migration/2 moved (direxec) to 0 CPU
migration/1 moved (direxec) to 2 CPU
migration/2 moved (direxec) to 1 CPU
migration/2 moved (direxec) to 0 CPU
migration/0 moved udiskd to 1 CPU
migration/2 moved bash to 1 CPU
migration/2 moved bash to 0 CPU
migration/1 moved (direxec) to 0 CPU
migration/1 moved (direxec) to 0 CPU
migration/2 moved (direxec) to 0 CPU
migration/1 moved (direxec) to 0 CPU
migration/1 moved (direxec) to 0 CPU
```

---

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

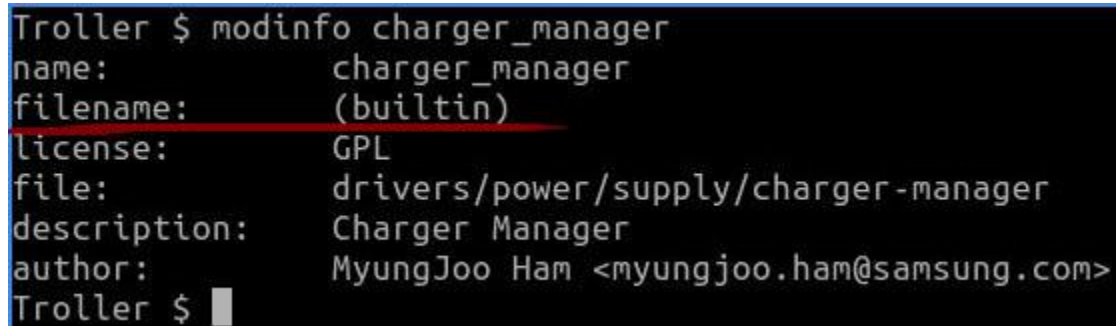
In summary, what the kernel thread “migration” does is to move threads from highly loaded vCPUs to others which are less crowded (by inserting them to a different run-queue). A function which is used by “migration” in order to move a task to a new run-queue is “move\_queued\_task” (<https://elixir.bootlin.com/linux/latest/source/kernel/sched/core.c#L2325>).

# charger-manager

The “charger\_manager” kernel thread is created by a freezable workqueue<sup>9</sup>. Freezable workqueues are basically frozen when the system is moved to a suspend state<sup>10</sup>. 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 memory<sup>11</sup>. The “Charger Manager” kernel module is written by MyungJoo Ham<sup>12</sup>.

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”<sup>13</sup>.

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 $
```

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

<sup>10</sup> <https://lwn.net/Articles/403891/>

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

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

<sup>13</sup> <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

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 workqueues in more details separately, for now it is all that you need to know).

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

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

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

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

The first one enables the tracing regarding workqueues. 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
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
k workerqueue=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
k workerqueue=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
k workerqueue=ata_sff req_cpu=8192 cpu=0
kworker/u2:1-48183 [000] d... 17860.134032: workqueue_queue_work: work struct=000000001e1ee94f function=kcryptd_crypt
dm_crypt] workqueue=kcryptd/253:0 req_cpu=8192 cpu=4294967295
kworker/u2:1-48183 [000] d... 17860.134074: workqueue_queue_work: work struct=00000000e0b6b12c function=kcryptd_crypt
dm_crypt] workqueue=kcryptd/253:0 req_cpu=8192 cpu=4294967295
```

# kdevtmpfs

“kdevtmpfs” is a kernel thread which was created using the “kthread\_run” function<sup>14</sup>. “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 filesystem<sup>15</sup>.

We can see the specific line of code that is used in order to create the mounting point “/dev”<sup>16</sup>. The mountpoint is created using the function “init\_mount”<sup>17</sup>. 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 code<sup>18</sup>.

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

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

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

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

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

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

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



## cpuhp

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"<sup>19</sup>.

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"<sup>20</sup> is called. In order to bring up a CPU function "cpu\_device\_up"<sup>21</sup> 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
```

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

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

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



# khungtaskd

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). The code of the kernel thread can be read in the following link [https://elixir.bootlin.com/linux/latest/source/kernel/hung\\_task.c](https://elixir.bootlin.com/linux/latest/source/kernel/hung_task.c).

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 outpted 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 called ([https://elixir.bootlin.com/linux/latest/source/kernel/hung\\_task.c#L178](https://elixir.bootlin.com/linux/latest/source/kernel/hung_task.c#L178)) - You can see the output in the screenshot below.

[illegible]

# 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 awakened by the page allocator, by calling “wakeup\_kswapd”<sup>22</sup>. The code of the page allocator is located at “mm/page\_alloc.c”<sup>23</sup>.

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()”<sup>24</sup>. In order to see that process happening we can use the following bpftrace one-liner: “**sudo bpftrace -e 'kfunc:balance\_pgdat { printf("%s:%d\n",curtask->comm,curtask->pid); }**” - You can see “kswapd0” calling it in the screenshot below. The flow of “kswapd” is based on limits, when to start shirking and “until when” to shrink (low and high limits).

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

---

<sup>22</sup> <https://elixir.bootlin.com/linux/latest/source/mm/vmscan.c#L4555>

<sup>23</sup> [https://elixir.bootlin.com/linux/latest/source/mm/page\\_alloc.c](https://elixir.bootlin.com/linux/latest/source/mm/page_alloc.c)

<sup>24</sup> <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 born<sup>25</sup>.

“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 migration<sup>26</sup> to do all the heavy lifting<sup>27</sup>. In order for “kcompactd” to work we should compile the kernel with “CONFIG\_COMPACTIOIN” enabled. Also, when a Linux system identifies that it is tight low in available memory the “kcompactd” won’t perform memory compaction memory<sup>28</sup>.

Overall, the “kcompactd” kernel thread is created in “kcompactd\_run” function<sup>29</sup> which is called by “kcompactd\_init”<sup>30</sup>. The function “kcompactd\_init” is started by “subsys\_initcall”<sup>31</sup>, which is responsible for initializing a subsystem.

The kernel thread starts the function “static int kcompactd(void \*p)”<sup>32</sup>. An instance of the kernel thread is created for each node (like vCPU) on the system<sup>33</sup>. 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 kswapd<sup>34</sup>. 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”<sup>35</sup> 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.

---

<sup>25</sup> <https://lwn.net/Articles/368869/>

<sup>26</sup> <https://lwn.net/Articles/157066/>

<sup>27</sup> <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L5>

<sup>28</sup> <https://www.linux-magazine.com/Issues/2015/179/Kernel-News>

<sup>29</sup> <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L2996>

<sup>30</sup> <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L3048>

<sup>31</sup> <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L3065>

<sup>32</sup> <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L2921>

<sup>33</sup> <https://elixir.bootlin.com/linux/v6.2-rc3/source/mm/compaction.c#L3061>

<sup>34</sup> <https://www.slideshare.net/AdrianHuang/memory-compaction-in-linux-kernelpdf>

<sup>35</sup> <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
```

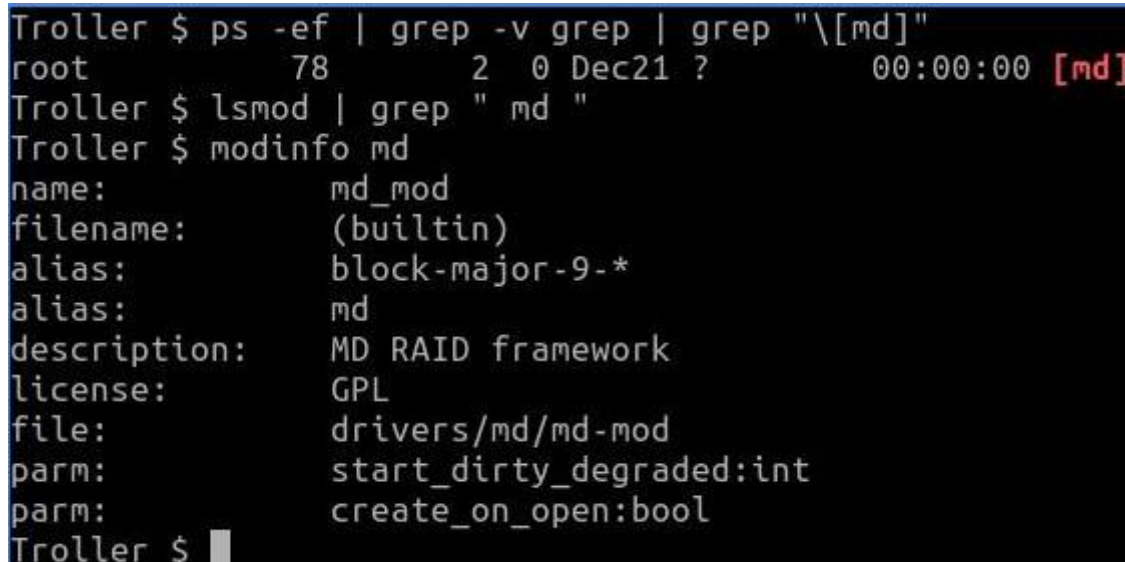
## md (Multiple Device Driver)

“md” is a kernel thread which is based on a workqueue<sup>36</sup>. 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 filesystem<sup>37</sup>.

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

Overall, “md” supports different RAID types: RAID 1 (mirroring), RAID 4, RAID 5, RAID 6 and RAID 10. For more information about RAID types I suggest reading the following link <https://www.prepressure.com/library/technology/raid>. Besides that, “md” also supports pseudo RAID technologies like: RAID 0, LINAR, MULTIPATH and FAULTY<sup>38</sup>.

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.

A terminal window with a black background and white text. The user 'Troller' runs several commands to check the md driver. The first command is 'ps -ef | grep -v grep | grep "\[md]"', which shows a process for 'md' with PID 78, PPID 2, and a status of '0'. The second command is 'lsmod | grep " md "', which shows the 'md' module is loaded. The third command is 'modinfo md', which displays detailed information about the md module, including its name, filename, alias, description, license, file path, and parameters.

```
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 $
```

<sup>36</sup> <https://elixir.bootlin.com/linux/v6.1/source/drivers/md/md.c#L9615>

<sup>37</sup> <https://linux.die.net/man/8/mdadm>

<sup>38</sup> [https://doxfer.webmin.com/Webmin/Linux\\_RAID](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])<sup>39</sup>. It can also be configured to allow access using “/dev/md/N” or “/dev/md/name”. If we want information about the current state of “md” we can query the file “/proc/mdstat” — for more information you can read <https://raid.wiki.kernel.org/index.php/Mdstat>. There is also the command line utility “mdadm” that can help with managing those devices<sup>40</sup>.

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

---

<sup>39</sup> <https://www.oreilly.com/library/view/managing-raid-on/9780596802035/ch01s03.html>

<sup>40</sup> <https://linux.die.net/man/8/mdadm>

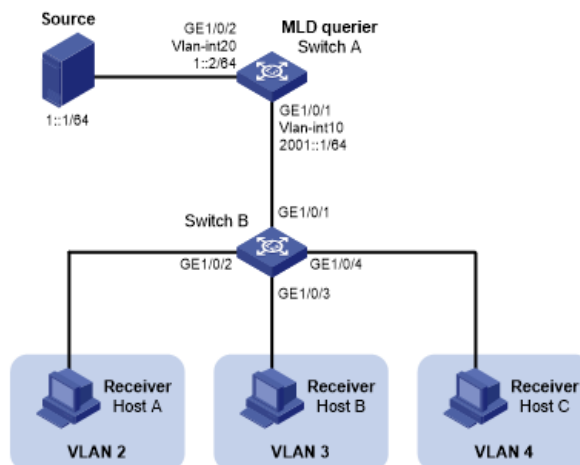
# mld (Multicast Listener Discovery)

“mld” is a kernel thread which was created using a workqueue<sup>41</sup>. 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 Windows<sup>42</sup> and Linux<sup>43</sup>.

We can think about it like IGMP<sup>44</sup> 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 types<sup>45</sup>.

Overall, MLD has three major message types: “Multicast Listener Query”, “Multicast Listener Report” and “Multicast Done”. For more information about them I suggest reading the following link<sup>46</sup>. Also, a more detailed explanation about the different MLD operations can be found in <https://ipccisco.com/lesson/mld-operations/>.

What “mld” does is to send MLD report messages<sup>47</sup> which are sent by an MLD host (see the diagram below<sup>48</sup>) and processes messages<sup>49</sup>. From the source code we can see that there are definitions for structs representing both MLDv1 and MLDv2 headers.



<sup>41</sup> <https://elixir.bootlin.com/linux/latest/source/net/ipv6/mcast.c#L3185>

<sup>42</sup> <https://learn.microsoft.com/en-us/windows/win32/winsock/igmp-and-windows-sockets>

<sup>43</sup> <https://lwn.net/Articles/29489/>

<sup>44</sup> <https://www.cloudflare.com/learning/network-layer/what-is-igmp/>

<sup>45</sup> <https://www.ibm.com/docs/en/zos/2.2.0?topic=protocol-multicast-listener-discovery>

<sup>46</sup> <https://community.cisco.com/t5/networking-knowledge-base/multicast-listener-discovery-mld/ta-p/3112082>

<sup>47</sup> <https://elixir.bootlin.com/linux/latest/source/net/ipv6/mcast.c#L3185>

<sup>48</sup> [https://techhub.hpe.com/eginfolib/networking/docs/switches/5130ei/5200-3944\\_ip-multi\\_cg/content/images/image33.png](https://techhub.hpe.com/eginfolib/networking/docs/switches/5130ei/5200-3944_ip-multi_cg/content/images/image33.png)

<sup>49</sup> <https://elixir.bootlin.com/linux/latest/source/net/ipv6/mcast.c#L1359>

## kmsd (Kernel Same Page Merging)

The kernel thread “ksm” is also known as “Kernel Same Page Merging” (and “kmsd” 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 “kmsd” to save memory due to de-duplication we should compile the kernel with “CONFIG\_KSM=y”. It is important to understand that the sharing of identical pages is done even if they are not shared by fork(). If you want to go over “kmsd” source code you can use the following link - <https://elixir.bootlin.com/linux/latest/source/mm/ksm.c>.

The way “kmsd” works is as follows. Scanning main memory for frames (“physical pages”) holding identical data and collects the virtual memory address that they are mapped. “kmsd” leaves one of those frames and remaps each duplicate one to point to the same frame. Lastly, “kmsd” 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 “kmsd” is created using the function `kthread_run`<sup>50</sup>. 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). “kmsd” 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.

“kmsd” 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 “kmsd” from running but keep the merged pages. “1” means run “kmsd”. “2” means stop “kmsd” from running and unmerge all currently merge pages (however leave the mergeable areas registered for next time).

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

---

<sup>50</sup> <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 objects<sup>51</sup>.

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

Moreover, “ttm\_swap” is a single threaded workqueue as seen in the Linux source code<sup>52</sup>.

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 driver<sup>53</sup>.

TTM is at the end a kernel module, you can find the source code and the Makefile in the kernel source tree<sup>54</sup>. Based on the module source code it is written by Thomas Hellstrom and Jerome Glisse<sup>55</sup>. Also, it is described as “TTM memory manager subsystem (for DRM device)”<sup>56</sup>. 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 you can see in the screenshot below.

---

<sup>51</sup> <https://docs.kernel.org/gpu/drm-mm.html>

<sup>52</sup> [https://elixir.bootlin.com/linux/v5.12.19/source/drivers/gpu/drm/ttm/ttm\\_memory.c#L424](https://elixir.bootlin.com/linux/v5.12.19/source/drivers/gpu/drm/ttm/ttm_memory.c#L424)

<sup>53</sup> <https://www.systutorials.com/docs/linux/man/7-drm-ttm/>

<sup>54</sup> <https://elixir.bootlin.com/linux/v6.1-rc2/source/drivers/gpu/drm/ttm>

<sup>55</sup> [https://elixir.bootlin.com/linux/v6.1-rc2/source/drivers/gpu/drm/ttm/ttm\\_module.c#L89](https://elixir.bootlin.com/linux/v6.1-rc2/source/drivers/gpu/drm/ttm/ttm_module.c#L89)

<sup>56</sup> [https://elixir.bootlin.com/linux/v6.1-rc2/source/drivers/gpu/drm/ttm/ttm\\_module.c#L89](https://elixir.bootlin.com/linux/v6.1-rc2/source/drivers/gpu/drm/ttm/ttm_module.c#L89)

```
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
```

# watchdogd

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 support<sup>57</sup>.

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 restart<sup>58</sup>.

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

The software part is being conducted by the “watchdogd” (the software watchdog daemon) which opens “/dev/watchdog” and writes to it in order to postpone the restart of the system by the hardware - for more information you can read <https://linux.die.net/man/8/watchdog>. Examples for different watchdog drives/modules for specific chips can be found in the source tree of linux here <https://elixir.bootlin.com/linux/v6.0.11/source/drivers/watchdog>. Some examples are apple\_wdt (Apple’s SOC), ath79\_wdt (Atheros AR71XX/AR724X/AR913X) and w83977f\_wdt (Winbond W83977F I/O Chip).

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

---

<sup>57</sup> <https://github.com/torvalds/linux/blob/master/Documentation/watchdog/watchdog-api.rst>

<sup>58</sup> <https://linuxhint.com/linux-kernel-watchdog-explained/>

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

Lastly, if you want to see the status of the watchdog you can use the command `"wdctl"`<sup>59</sup> - As can be seen in the screenshot below<sup>60</sup>. For more information about the concept I suggest reading [https://en.wikipedia.org/wiki/Watchdog\\_timer](https://en.wikipedia.org/wiki/Watchdog_timer).

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

---

<sup>59</sup> <https://man7.org/linux/man-pages/man8/wdctl.8.html>

<sup>60</sup> [https://en.wikipedia.org/wiki/Watchdog\\_timer#/media/File:Wdctl\\_screenshot.png](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\swapfile<sup>61</sup>. 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 pool<sup>62</sup>.

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 cache<sup>63</sup>. The “zswap\_frontswap\_store” is the function that attempts to compress and store a single page<sup>64</sup>.

The kernel thread “zswap-shrink” is created based on a workqueue<sup>65</sup>. On my VM (Ubuntu 22.04.1) zswap is compiled part of the kernel itself and not as a separate “\*.ko” (kernel module). You can see in the screenshot below that it does not appear in the output of “lsmod” and is marked as builtin (look at the filename field) in the output of “modinfo”.

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

For more information like the compression used by zswap (the default one is lzo) and other parameters that can be configured for zswap I suggest reading the following link

---

<sup>61</sup> <https://www.kernel.org/doc/html/v4.18/vm/frontswap.html>

<sup>62</sup> <https://elixir.bootlin.com/linux/latest/source/mm/zswap.c>

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

<sup>64</sup> <https://elixir.bootlin.com/linux/v6.1-rc2/source/mm/zswap.c#L1097>

<sup>65</sup> <https://elixir.bootlin.com/linux/v6.1-rc2/source/mm/zswap.c#L1511>

<https://wiki.archlinux.org/title/zswap>. You can also read the parameters “/sys/module/zswap/parameters”.

## khugepaged

The kernel thread “khugepaged” is created using the “kthread\_run()” function<sup>66</sup>. It is responsible for the “Transparent Hugepage Support” (aka THP). “khugepaged” scans memory and collapses sequences of basic pages into huge pages<sup>67</sup>.

We can manage and configure TPH using sysfs<sup>68</sup> or by using the syscalls “madvise”<sup>69</sup> and “prctl”<sup>70</sup>. The scan of memory is done by calling “khugepaged\_do\_scan()”<sup>71</sup> which in turn calls “khugepaged\_scan\_mm\_slot()”<sup>72</sup>. In order to demonstrate that I have used the following bpftrace oneliner **“sudo bpftrace -e 'kfunc:khugepaged\_scan\_mm\_slot{ printf("%s:%d\n",curtask->comm,curtask->pid); }”**. The output is shown in the screenshot below.

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

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

---

<sup>66</sup> <https://elixir.bootlin.com/linux/latest/source/mm/khugepaged.c#L2551>

<sup>67</sup> <https://www.kernel.org/doc/html/latest/admin-guide/mm/transhuge.html>

<sup>68</sup> <https://www.kernel.org/doc/html/latest/admin-guide/mm/transhuge.html#thp-sysfs>

<sup>69</sup> <https://man7.org/linux/man-pages/man2/madvise.2.html>

<sup>70</sup> <https://man7.org/linux/man-pages/man2/prctl.2.html>

<sup>71</sup> <https://elixir.bootlin.com/linux/latest/source/mm/khugepaged.c#L2404>

<sup>72</sup> <https://elixir.bootlin.com/linux/v6.1.12/source/mm/khugepaged.c#L2250>

# krfcommd

“krfcommd” is a kernel which is started by executing “kthread\_run()” function<sup>73</sup>. The kernel thread executes the “rfcomm\_run()” function<sup>74</sup>. Thus, we can say that “krfcommd” is responsible for RFCOMM connections<sup>75</sup>.

RFCOMM (Radio Frequency Communication) is a set of transport protocols on top of L2CAP which provides emulated RS-232 serial ports. It provides a simple reliable data stream (like TCP). It is used directly by many telephony related profiles as a carrier for AT commands, as well as being a transport layer for OBEX over Bluetooth<sup>76</sup>.

Moreover, there is also an “rfcomm” cli tool in Linux. It is used to inspect and maintain RFCOMM configuration<sup>77</sup>. For more information about RFCOMM I suggest reading <https://www.btframework.com/rfcomm.htm>. You can also go over the protocol specification<sup>78</sup>.

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 them<sup>79</sup>.

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:     2787EECAEC282A1A24A7701
depends:         bluetooth
retpoline:      Y
intree:         Y
name:           rfcomm
vermagic:       5.19.7-arch1-1.0 SMP preempt mod_unload 686
sig_id:         PKCS#7
signer:         Build time autogenerated kernel key
sig_key:        30:9A:19:01:BA:9C:BA:D5:C0:8D:F7:A5:39:AA:C7:54:A6:C9:D8:2B
sig_hashalgo:   sha512
signature:      30:64:02:30:6C:AB:DA:07:56:CC:36:9D:66:06:E2:8B:98:E9:4A:50:
77:C0:37:08:0A:12:CD:5D:84:F7:2F:4A:FA:CB:58:68:B9:C4:7B:C0:
08:1C:EC:61:33:FA:7E:A8:69:6B:FD:E7:02:30:69:C8:06:98:12:9C:
E3:B3:25:33:03:12:81:D6:77:59:54:F5:8E:5B:D5:FF:C4:5D:D1:F1:
02:0E:16:68:2E: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)
```

<sup>73</sup> <https://elixir.bootlin.com/linux/latest/source/net/bluetooth/rfcomm/core.c#L2215>

<sup>74</sup> <https://elixir.bootlin.com/linux/latest/source/net/bluetooth/rfcomm/core.c#L2109>

<sup>75</sup> <https://stackoverflow.com/questions/57152408/what-is-the-internal-mechanics-of-socket-function>

<sup>76</sup> [https://en.wikipedia.org/wiki/List\\_of\\_Bluetooth\\_protocols](https://en.wikipedia.org/wiki/List_of_Bluetooth_protocols)

<sup>77</sup> <https://linux.die.net/man/1/rfcomm>

<sup>78</sup> <https://www.bluetooth.com/specifications/specs/rfcomm-1-1/>

<sup>79</sup> [https://www.amd.e-technik.uni-rostock.de/ma/gol/lectures/wirlec/bluetooth\\_info/rfcomm.html](https://www.amd.e-technik.uni-rostock.de/ma/gol/lectures/wirlec/bluetooth_info/rfcomm.html)



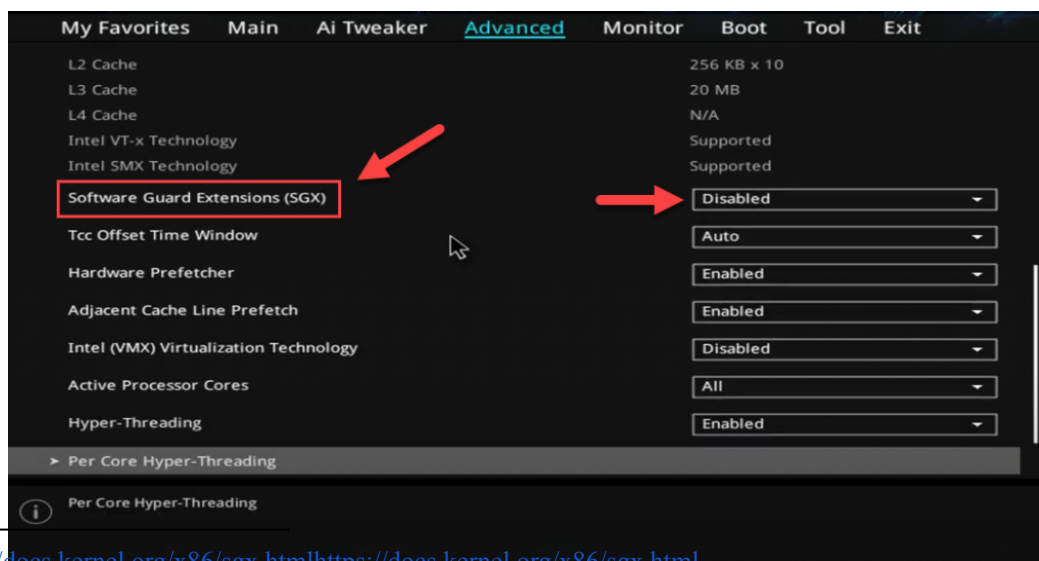
# ksgxd

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 regions<sup>80</sup>. For more information about SGX you can read my writeup about it<sup>81</sup>.

“ksgxd” is a kernel which is started by executing “kthread\_run()” function<sup>82</sup>. The kernel thread executes the “ksgxd” function<sup>83</sup>. “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 memory<sup>84</sup> like “kswapd”<sup>85</sup>.

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 below<sup>86</sup>.

Lastly, there is a great guide for an example SGX app using a Linux VM on Azure that I encourage you to read<sup>87</sup>. For more information about the Linux stack for SGX I suggest reading [https://download.01.org/intelsgxstack/2021-12-08/Getting\\_Started.pdf](https://download.01.org/intelsgxstack/2021-12-08/Getting_Started.pdf) and going over the following github repo <https://github.com/intel/linux-sgx>.



<sup>80</sup> <https://docs.kernel.org/x86/sgx.html>

<sup>81</sup> <https://medium.com/@boutnaru/security-sgx-software-guard-extension-695cab7dbcb2>

<sup>82</sup> <https://elixir.bootlin.com/linux/v6.1.10/source/arch/x86/kernel/cpu/sgx/main.c#L427>

<sup>83</sup> <https://elixir.bootlin.com/linux/v6.1.10/source/arch/x86/kernel/cpu/sgx/main.c#L395>

<sup>84</sup> <https://elixir.bootlin.com/linux/v6.1.10/source/arch/x86/kernel/cpu/sgx/main.c#L188>

<sup>85</sup> <https://medium.com/@boutnaru/the-linux-process-journey-kswapd-22754e783901>

<sup>86</sup> <https://phoenixnap.com/kb/intel-sgx>

<sup>87</sup> <https://tsmatz.wordpress.com/2022/05/17/confidential-computing-intel-sgx-enclave-getting-started/>

## jbd2 (Journal Block Device 2)

“JBD” stands for “Journal Block Device”<sup>88</sup>. “jbd2” is a kernel which is started by executing “kthread\_run()” function<sup>89</sup>. 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 log<sup>90</sup>.

The kernel thread executes the “kjournald2()” function<sup>91</sup>. 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 file<sup>92</sup>.

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 ext4<sup>93</sup>.

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

<sup>88</sup> <https://manpages.ubuntu.com/manpages/jammy/man1/pmdajbd2.1.html>

<sup>89</sup> <https://elixir.bootlin.com/linux/v6.2.1/source/fs/jbd2/journal.c#L277>

<sup>90</sup> <https://elixir.bootlin.com/linux/v6.2.1/source/fs/jbd2/journal.c>

<sup>91</sup> <https://elixir.bootlin.com/linux/v6.2.1/source/fs/jbd2/journal.c#L169>

<sup>92</sup> <https://elixir.bootlin.com/linux/v6.2.1/source/fs/jbd2/journal.c#L152>

<sup>93</sup> [https://en.wikipedia.org/wiki/Journaling\\_block\\_device](https://en.wikipedia.org/wiki/Journaling_block_device)

# netns

The kernel thread “netns” is based on a single threaded workqueue<sup>94</sup>, which is created when the network namespace is initialized (net\_ns\_init()). If you want to read more about “network namespaces” you can use the following link

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

Also, for a reminder you can also check out the diagram below<sup>95</sup>.

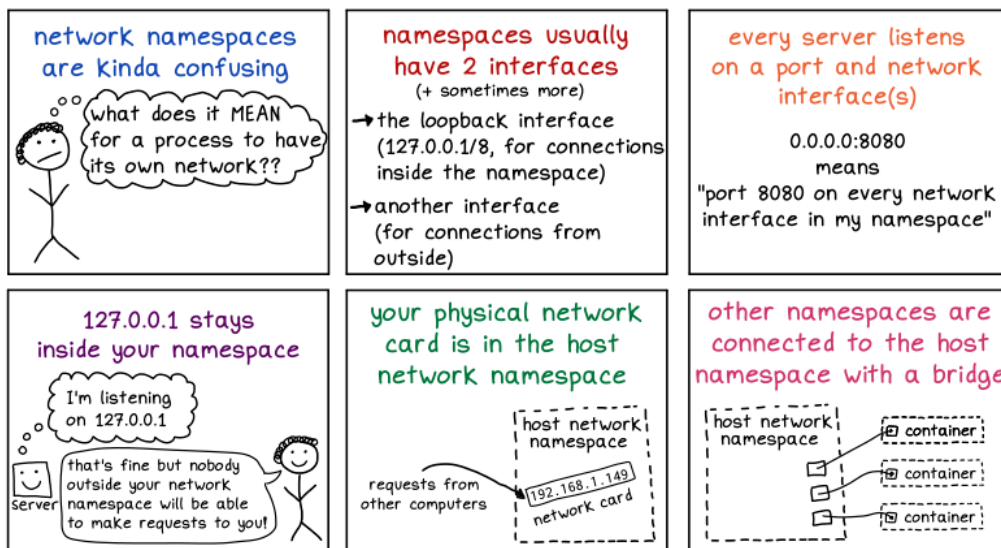
“netns” is responsible for cleaning up network namespaces. When a namespace is destroyed the kernel adds it to a cleanup list. The kernel thread “netns” goes over the list and performs the cleanup process using the “cleanup\_net()” function<sup>96</sup>.

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()” function<sup>97</sup>.

SULIA EVANS  
@b0rk

## network namespaces

18



<sup>94</sup> [https://elixir.bootlin.com/linux/v6.2-rc4/source/net/core/net\\_namespace.c#L1106](https://elixir.bootlin.com/linux/v6.2-rc4/source/net/core/net_namespace.c#L1106)

<sup>95</sup> <https://wizardzines.com/comics/network-namespaces/>

<sup>96</sup> [https://elixir.bootlin.com/linux/v6.2.3/source/net/core/net\\_namespace.c#L565](https://elixir.bootlin.com/linux/v6.2.3/source/net/core/net_namespace.c#L565)

<sup>97</sup> [https://elixir.bootlin.com/linux/v6.2-rc4/source/net/core/net\\_namespace.c#L649](https://elixir.bootlin.com/linux/v6.2-rc4/source/net/core/net_namespace.c#L649)

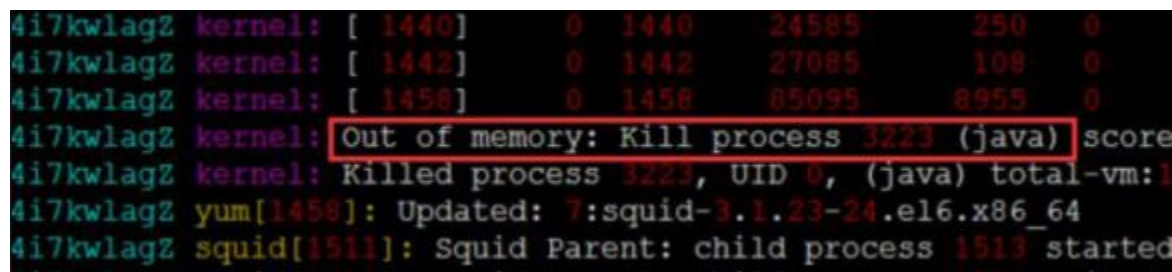
## oom\_reaper

“oom\_reaper” is a kernel thread which was created using the “kthread\_run” function<sup>98</sup>. Basically, it is the implementation of the OMM (Out-of-Memory) killer function of the Linux kernel - for more information about it I encourage you to read the following link <https://medium.com/@boutnaru/linux-out-of-memory-killer-oom-killer-bb2523da15fc>.

The function which is executed by the thread is “oom\_reaper”<sup>99</sup> which calls “oom\_reap\_task”<sup>100</sup>.

Based on the documentation the goal of the “oom\_reaper” kernel thread is to try and reap the memory used by the OOM victim<sup>101</sup>. “oom\_reaper” sleeps until it is waked up<sup>102</sup> which is after OOM kills the process<sup>103</sup>.

After killing the process the victim is queued so the “oom\_reaper” can release the resources<sup>104</sup>. You can see an example of the log created by OOM after killing a process<sup>105</sup>.



```
4i7kwlagZ kernel: [ 1440] 0 1440 24585 250 0
4i7kwlagZ kernel: [ 1442] 0 1442 27085 108 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
```

<sup>98</sup> [https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom\\_kill.c#L735](https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L735)

<sup>99</sup> [https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom\\_kill.c#L640](https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L640)

<sup>100</sup> [https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom\\_kill.c#L609](https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L609)

<sup>101</sup> [https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom\\_kill.c#L504](https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L504)

<sup>102</sup> [https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom\\_kill.c#L680](https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L680)

<sup>103</sup> [https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom\\_kill.c#L947](https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L947)

<sup>104</sup> [https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom\\_kill.c#L992](https://elixir.bootlin.com/linux/v6.2.5/source/mm/oom_kill.c#L992)

<sup>105</sup> <https://blog.capdata.fr/index.php/linux-out-of-memory-killer-oom-killer-pour-un-serveur-base-de-donnees-postgresql/>

# kpsmoused

“kpsmoused” is a kernel thread which based on an ordered workqueue<sup>106</sup> which is allocated inside the “pmouse\_init” function. “kpsmoused” is responsible for handling the input from PS/2 mouse devices.

Thus, “kpsmoused” transforms the raw data to high level event of mouse movements that be can consume from “/dev/input/mice”, “/dev/input/mouseX”, or “/dev/input/eventX”<sup>107</sup>.

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 2.5.72<sup>108</sup>.

```
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:        Vojtech Pavlik <vojtech@suse.cz>
srcversion:    900B3C75C18D7DDE1706120
alias:         serio:ty05pr=id*ex*
alias:         serio:ty01pr=id*ex*
depends:        libps2,serio
retpoline:     y
intree:        y
name:          psmouse
vermagic:      5.19.7-arch1-1.0 SMP preempt mod_unload 686
sig_id:        PKCS#7
signer:        Build time autogenerated kernel key
sig_key:       30:9a:19:01:ba:9c:ba:d5:c0:8d:f7:a5:39:aa:c7:54:a6:c9:d0:2b
sig_hashalgo:  sha512
signature:     30:64:02:30:26:9e:10:64:df:8e:1f:2e:c6:2d:a8:f3:e1:53:a5:4a:
8f:01:f0:6b:f7:e7:a1:02:f4:6a:48:d0:43:fb:7b:3a:7e:0a:25:7b:
35:85:6d:cc:22:2c:89:4a:93:2a:fc:b0:02:30:0b:da:0b:5e:c8:7d:
bd:96:ae:66:72:d4:4c:a7:64:59:e5:1d:76:0b:04:f2:20:70:93:d0:
23:ff:8b:9f:57:f5:d2:0b:9c:98:ea:37:08:d5:39:72:3b:50:87:7c:
13:c5
parm:          tpdebug:enable debugging, dumping packets to KERN_DEBUG. (bool)
parm:          recalib_delta:packets containing a delta this large will be discarded, and a recalibration may be scheduled. (int)
parm:          jumpy_delay:delay (ms) before recal after jumpiness detected (int)
parm:          spew_delay:delay (ms) before recal after packet spew detected (int)
parm:          recal_guard_time:interval (ms) during which recal will be restarted if packet received (int)
parm:          post_interrupt_delay:delay (ms) before recal after recal interrupt detected (int)
parm:          autorecal:enable recalibration in the driver (bool)
parm:          hoggk_mode:default hoggk 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:          atech_workaround:AtTech second scroll wheel workaround, 1 = enabled, 0 = disabled (default). (bool)
parm:          resetafter:Reset device after so many bad packets (0 = never). (uint)
parm:          resync_time:How long can mouse stay idle before forcing resync (in seconds, 0 = never). (uint)
```

<sup>106</sup> <https://elixir.bootlin.com/linux/v6.2.6/source/drivers/input/mouse/psmouse-base.c#L2046>

<sup>107</sup> <https://www.kernel.org/doc/html/v5.5/input/input.html>

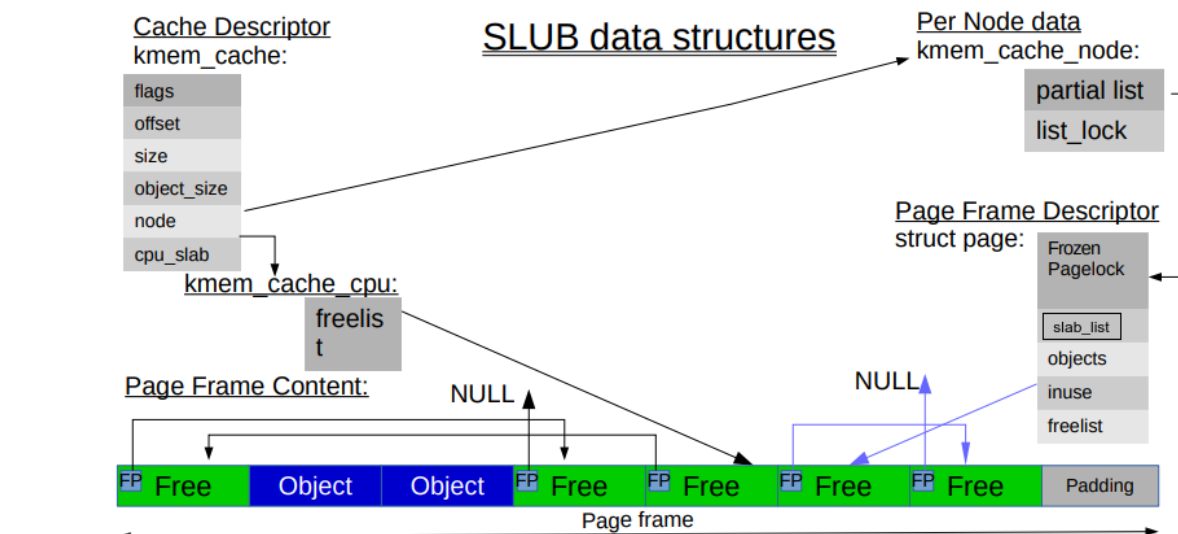
<sup>108</sup> <https://elixir.bootlin.com/linux/v2.5.72/source/drivers/input/mouse/psmouse-base.c>

## slub\_flushwq

“slub\_flushwq” is a kernel thread which based on a workqueue<sup>109</sup> 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 enabled<sup>110</sup>. 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 RAM<sup>111</sup>. The queuing of work is done inside the “flush\_all\_cpus\_locked” function<sup>112</sup>.

SLUB is also known as the “Unqueued Slab Allocator”<sup>113</sup>. Slab allocation is a memory management mechanism which allows efficient memory allocation of objects. It is done using reduction of fragmentation that is caused due to allocations/deallocations<sup>114</sup>. For more information about slab allocation I suggest reading the following link <https://hammertux.github.io/slab-allocator>.

Thus, SLUB is a slab allocator that limits the use of cache lines instead of using queued object per cpu/per node list<sup>115</sup>. So, it is less complicated because it does not keep queues (like for each CPU). The only queue is a linked list for all the objects in each of the slub pages<sup>116</sup>. 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 below<sup>117</sup>.



<sup>109</sup> <https://elixir.bootlin.com/linux/v6.2.6/source/mm/slub.c#L5057>

<sup>110</sup> <https://elixir.bootlin.com/linux/v6.2.6/source/mm/slub.c#L5056>

<sup>111</sup> [https://cateee.net/lkddb/web-lkddb/SLUB\\_TINY.html](https://cateee.net/lkddb/web-lkddb/SLUB_TINY.html)

<sup>112</sup> <https://elixir.bootlin.com/linux/v6.2.6/source/mm/slub.c#L2822>

<sup>113</sup> <https://lwn.net/Articles/229096/>

<sup>114</sup> [https://en.wikipedia.org/wiki/Slab\\_allocation](https://en.wikipedia.org/wiki/Slab_allocation)

<sup>115</sup> <https://elixir.bootlin.com/linux/v6.2.6/source/mm/slub.c#L3>

<sup>116</sup> <https://hammertux.github.io/slab-allocator>

<sup>117</sup> <https://hammertux.github.io/img/SLUB-DS.png>

# pgdatinit

“pgdatinit” is a kernel which is started by executing the “kthread\_run()” function<sup>118</sup>. The kernel thread executes the “deferred\_init\_memmap()” function<sup>119</sup>.

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]”<sup>120</sup>.

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 threads<sup>121</sup>.

Lastly, after the initialization flow is finished an information message is sent to the kernel ring buffer<sup>122</sup> - as you can see in the image below<sup>123</sup>.

```
[ 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
```

---

<sup>118</sup> [https://elixir.bootlin.com/linux/v6.3-rc4/source/mm/page\\_alloc.c#L2284](https://elixir.bootlin.com/linux/v6.3-rc4/source/mm/page_alloc.c#L2284)

<sup>119</sup> [https://elixir.bootlin.com/linux/v6.3-rc4/source/mm/page\\_alloc.c#L2108](https://elixir.bootlin.com/linux/v6.3-rc4/source/mm/page_alloc.c#L2108)

<sup>120</sup> [https://elixir.bootlin.com/linux/v6.3-rc4/source/mm/page\\_alloc.c#L2283](https://elixir.bootlin.com/linux/v6.3-rc4/source/mm/page_alloc.c#L2283)

<sup>121</sup> [https://cateee.net/lkddb/web-lkddb/DEFERRED\\_STRUCT\\_PAGE\\_INIT.html](https://cateee.net/lkddb/web-lkddb/DEFERRED_STRUCT_PAGE_INIT.html)

<sup>122</sup> [https://elixir.bootlin.com/linux/v6.3-rc4/source/mm/page\\_alloc.c#L2177](https://elixir.bootlin.com/linux/v6.3-rc4/source/mm/page_alloc.c#L2177)

<sup>123</sup> <https://www.mail-archive.com/debian-bugs-dist@lists.debian.org/msg1822096.html>

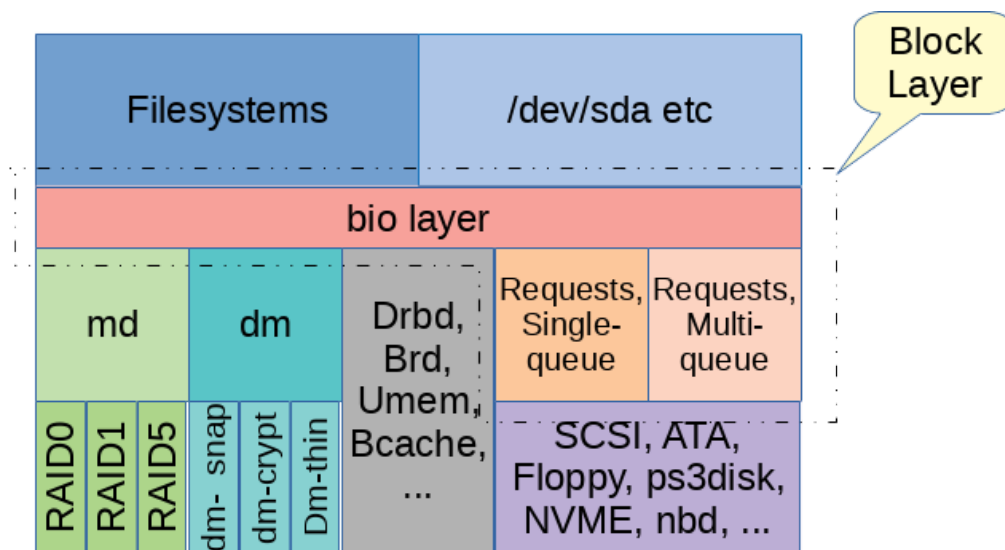
# kblockd

“kblockd” is a kernel thread based on a workqueue<sup>124</sup> 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 below<sup>125</sup>.

Overall, one might think that we can use keventd<sup>126</sup> 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 drivers<sup>127</sup>.

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 devices<sup>128</sup>.



<sup>124</sup> <https://elixir.bootlin.com/linux/v6.2.9/source/block/blk-core.c#L1191>

<sup>125</sup> <https://lwn.net/Articles/736534/>

<sup>126</sup> <https://lwn.net/Articles/11351/>

<sup>127</sup> <https://mirrors.edge.kernel.org/pub/linux/kernel/people/akpm/patches/2.5/2.5.70/2.5.70-mm8/broken-out/kblockd.patch>

<sup>128</sup> <https://elixir.bootlin.com/linux/v6.3-rc4/source/block/blk-core.c#L13>

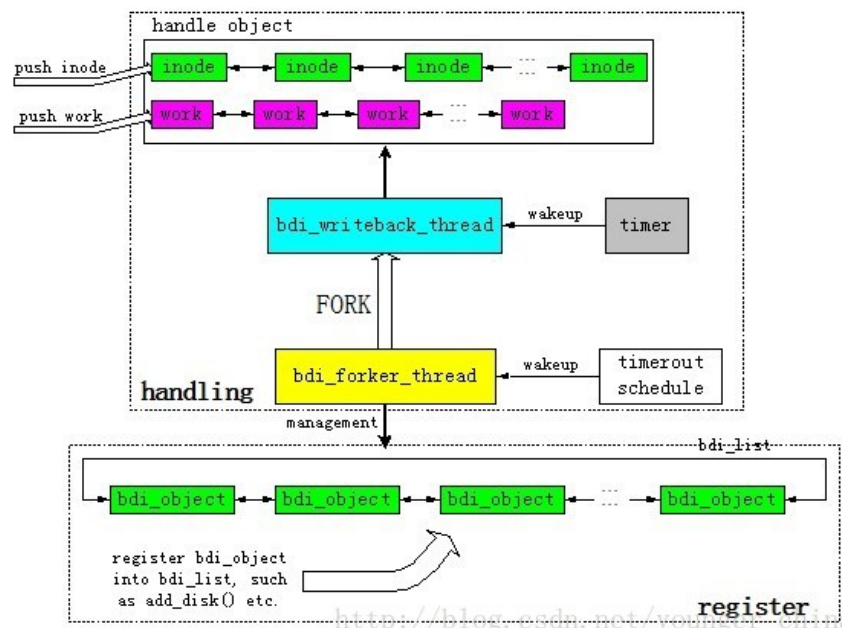


# writeback

The kernel thread “writeback” is based on a workqueue<sup>129</sup>. The goal of the kernel thread is to serve all async writeback tasks<sup>130</sup>. 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 disk<sup>131</sup>.

Overall, they are two ways of flushing dirty pages using writeback. The first is in case of an explicit writeback request - like syncing inode pages of a superblock. Thus, the “wb\_start\_writeback()” is called with the superblock information and the number of pages to flush. The second one is when there is no specific writeback request, in this case there is a timer that wakes up the thread periodically to flush dirty data<sup>132</sup>.

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 disk<sup>133</sup>. By the way, “bdi” stands for “Backing Device Information”<sup>134</sup>. Lastly, to get an overview of the “writeback” mechanism you can checkout the diagram below<sup>135</sup>.



<sup>129</sup> <https://elixir.bootlin.com/linux/v6.2.5/source/mm/backing-dev.c#L303>

<sup>130</sup> <https://elixir.bootlin.com/linux/v6.2.5/source/mm/backing-dev.c#L35>

<sup>131</sup> <https://www.oreilly.com/library/view/understanding-the-linux/0596005652/ch15s01.html>

<sup>132</sup> <https://lwn.net/Articles/326552/>

<sup>133</sup> [https://blog.csdn.net/younger\\_china/article/details/55187057](https://blog.csdn.net/younger_china/article/details/55187057)

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

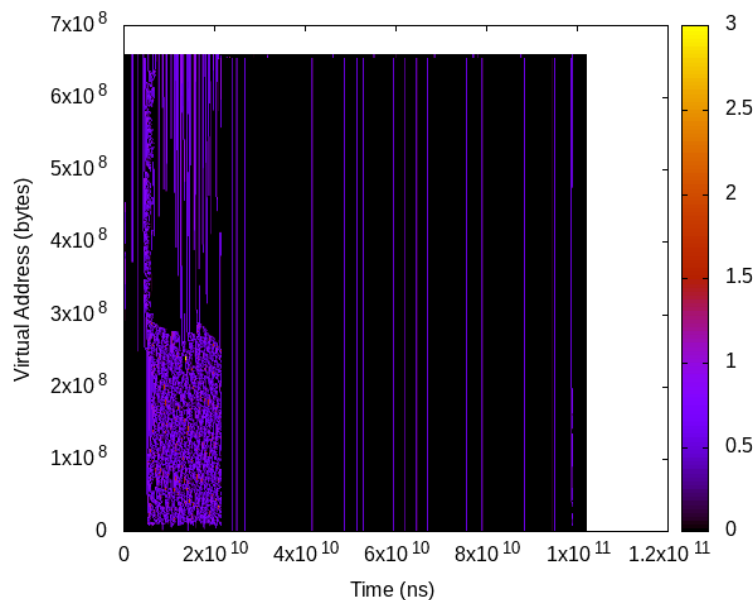
<sup>135</sup> [https://blog.csdn.net/younger\\_china/article/details/55187057](https://blog.csdn.net/younger_china/article/details/55187057)

## kdamond (Data Access MONitor)

“kdamond” is a kernel thread which is created using the “kthread\_run()” function<sup>136</sup> which is part of the DAMON (Data Access MONitor) subsystem. The kernel thread executes the “kdamon\_fn()” function<sup>137</sup>. Overall, DAMON provides a lightweight data access monitoring facility that can help users in analyzing the memory access patterns of their systems<sup>138</sup>. Based on the documentation DAMON increases the memory usage by 0.12% and slows the workloads down by 1.39%<sup>139</sup>.

Also, DAMON has an API for kernel programs<sup>140</sup>. 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 configurations<sup>141</sup>.

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 benchmarks<sup>142</sup>. 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-blackscholes<sup>143</sup>. Lastly, there are also other mechanisms in Linux that can help with data access monitoring such as “Perf Mem” and “Idle Page Tracking”



<sup>136</sup> <https://elixir.bootlin.com/linux/v6.3-rc5/source/mm/daemon/core.c#L632>

<sup>137</sup> <https://elixir.bootlin.com/linux/v6.3-rc5/source/mm/daemon/core.c#L1304>

<sup>138</sup> <https://www.kernel.org/doc/html/latest/admin-guide/mm/daemon/index.html>

<sup>139</sup> <https://damonitor.github.io/doc/html/v20/vm/daemon/eval.html>

<sup>140</sup> <https://www.kernel.org/doc/html/v5.17/vm/daemon/api.html#functions>

<sup>141</sup> <https://sjp38.github.io/post/daemon/>

<sup>142</sup> <https://parsec.cs.princeton.edu/parsec3-doc.htm>

<sup>143</sup> <https://lwn.net/Articles/813108/>

# kintegrityd

“kintegrityd” is a kernel thread based on a workqueue<sup>144</sup> 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”<sup>145</sup>. The function is called to complete a read request by verifying the transferred integrity metadata and then calls the original bio end\_io function<sup>146</sup>.

This procedure is done to ensure that the data was not changed by mistake (like in a case of a bug or an hardware failure<sup>147</sup>. 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 layer<sup>148</sup>.

Moreover, in order for it to work we should enable CONFIG\_BLK\_DEV\_INTEGRITY, which is defined as “Block layer data integrity support”<sup>149</sup>. 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() function<sup>150</sup>. We can toggle the writing of metadata using “/sys/block/<BlockDevice>/integrity/write\_generate“ and the verification of the metadata using “/sys/block/<BlockDevice>/integrity/read\_verify” - as shown in the screenshot below.

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

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

<sup>144</sup> <https://elixir.bootlin.com/linux/v6.1/source/block/bio-integrity.c#L455>

<sup>145</sup> <https://elixir.bootlin.com/linux/v6.1/source/block/bio-integrity.c#L317>

<sup>146</sup> <https://elixir.bootlin.com/linux/v6.1/source/block/bio-integrity.c#L313>

<sup>147</sup> <https://www.quora.com/What-is-the-purpose-of-kintegrityd-Linux-Kernel-Daemon/answer/Liran-Ben-Haim>

<sup>148</sup> <https://www.kernel.org/doc/Documentation/block/data-integrity.txt>

<sup>149</sup> <https://elixir.bootlin.com/linux/v6.1/source/block/Kconfig#L60>

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

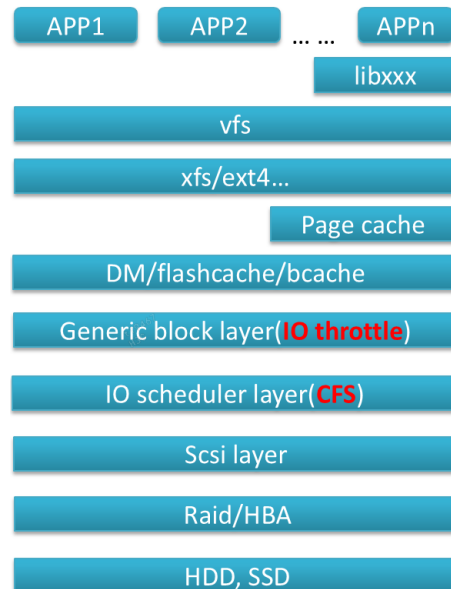
# kthrotld

“kthrotld” is a kernel thread which was created using an workqueue<sup>151</sup> 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 queues<sup>152</sup>.

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”<sup>153</sup> and then where is it referenced<sup>154</sup>. 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 device<sup>155</sup>.

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)<sup>156</sup>.

Overall, IO throttling is done as part of the generic block layer and before the IO scheduler as seen in the diagram below<sup>157</sup>. For more information on “Block Throttling” I suggest reading <https://developer.aliyun.com/article/789736>.



<sup>151</sup> <https://elixir.bootlin.com/linux/v6.1.1/source/block/blk-throttle.c#L2470>

<sup>152</sup> <https://www.halolinux.us/kernel-architecture/request-queues.html>

<sup>153</sup> <https://elixir.bootlin.com/linux/v6.1.1/source/include/linux/blkdev.h#L395>

<sup>154</sup> [https://elixir.bootlin.com/linux/v6.1.1/C/ident/request\\_queue](https://elixir.bootlin.com/linux/v6.1.1/C/ident/request_queue)

<sup>155</sup> <https://www.oreilly.com/library/view/linux-device-drivers/0596000081/ch12s04.html>

<sup>156</sup> <https://developer.aliyun.com/article/789736>

<sup>157</sup> <https://blog.csdn.net/yiyeguzhou100/article/details/104044419>

## scsi\_eh (Small Computer System Interface Error Handling)

The kernel thread “scsi\_eh” is executed using the “kthread\_run” function. The name pattern of the kernel thread is “scsi\_eh\_<SCSI\_HOST\_NUMBER>”<sup>158</sup>. It is the “SCSI error handler” which is responsible for all of the error handling targeting every SCSI host<sup>159</sup>. The kernel thread is executing the “scsi\_error\_handler” function<sup>160</sup>.

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 below<sup>161</sup>.

Lastly, SCSI stands for “Small Computer System Interface”. It is a set of standards (from ANSI) for electronic interfaces in order to communicate with peripheral hardware like CD-ROM drives, tape drives, printers, disk drives and more<sup>162</sup>.. For more information about SCSI I suggest going over <https://hackaday.com/2023/03/02/scsi-the-disk-bus-for-everything/>.



---

<sup>158</sup> <https://elixir.bootlin.com/linux/v6.4-rc1/source/drivers/scsi/hosts.c#L504>

<sup>159</sup> [https://elixir.bootlin.com/linux/v6.4-rc1/source/drivers/scsi/scsi\\_error.c#L2230](https://elixir.bootlin.com/linux/v6.4-rc1/source/drivers/scsi/scsi_error.c#L2230)

<sup>160</sup> [https://elixir.bootlin.com/linux/v6.4-rc1/source/drivers/scsi/scsi\\_error.c#L2233](https://elixir.bootlin.com/linux/v6.4-rc1/source/drivers/scsi/scsi_error.c#L2233)

<sup>161</sup> <https://computer.howstuffworks.com/scsi.htm>

<sup>162</sup> <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” function<sup>163</sup>. It is part of the common block controller cgroup interface<sup>164</sup>.

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 cgroup<sup>165</sup>. An example of priority inversion is shown in the diagram below<sup>166</sup>.

Thus, to avoid the problem mentioned above the function “submit\_bio”<sup>167</sup> punts the issuing of the bio request to a dedicated work item (per-block cgroup).

It calls “blkcg\_punt\_bio\_submit”<sup>168</sup>, which will call “\_\_blkcg\_punt\_bio\_submit”<sup>169</sup>.

### Priority inversion.

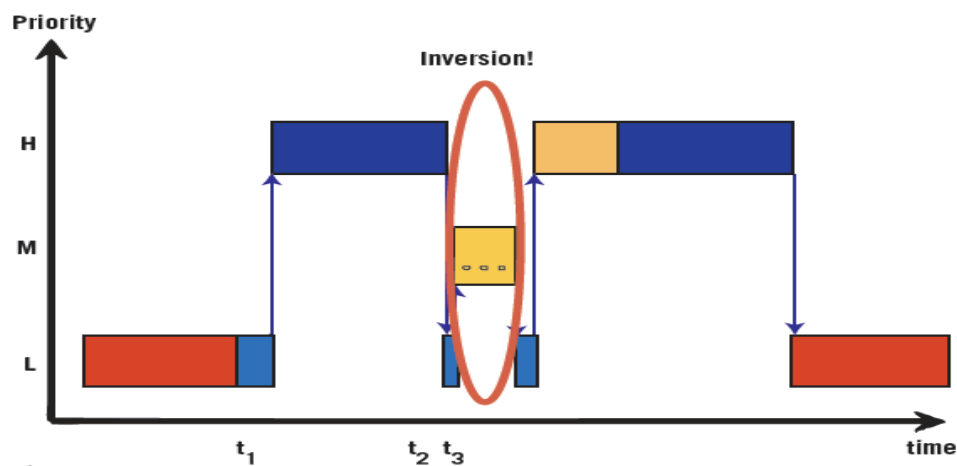


Figure 2

<sup>163</sup> <https://elixir.bootlin.com/linux/v6.2.5/source/block/blk-cgroup.c#L2058>

<sup>164</sup> <https://elixir.bootlin.com/linux/v6.2.5/source/block/blk-cgroup.c#L3>

<sup>165</sup> <https://patchwork.kernel.org/project/linux-block/patch/20190627203952.386785-6-tj@kernel.org/>

<sup>166</sup> <https://embeddedgurus.com/barr-code/2010/11/firmware-specific-bug-8-priority-inversion/>

<sup>167</sup> <https://elixir.bootlin.com/linux/v6.2.5/source/block/blk-core.c#L829>

<sup>168</sup> <https://elixir.bootlin.com/linux/v6.2.5/source/block/blk-cgroup.h#L380>

<sup>169</sup> <https://elixir.bootlin.com/linux/v6.2.5/source/block/blk-cgroup.c#L1657>