

# **The Linux Kernel Data Structures Journey**

**version 1.0 (beta)  
August-2023**

**By Shlomi Boutnaru**



Created using [Craiyon AI Image Generator](#)

<b>Introduction</b>	<b>3</b>
<b>struct task_struct</b>	<b>4</b>
<b>struct mm_struct</b>	<b>5</b>
<b>struct vm_area_struct</b>	<b>6</b>
<b>struct therad_info</b>	<b>7</b>
<b>struct thread_struct</b>	<b>8</b>
<b>struct inode</b>	<b>9</b>

# Introduction

When starting to read the source code of the Linux kernel I believe that they are basic data structures that everyone needs to know about. Because of that I have decided to write a series of short writeups aimed at providing the basic vocabulary and understanding for achieving that.

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

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!!!!!!

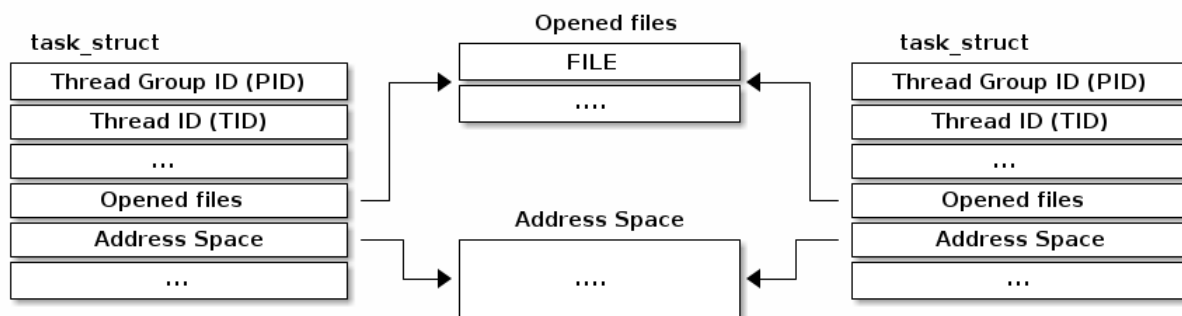
# struct task\_struct

Every operating system has a data structure that represents a “process<sup>1</sup> object” (generally called PCB - Process Control Block). By the way, “task\_struct” is the PCB in Linux (it is also the TCB, meaning the Thread Control Block). As an example, a diagram that shows two processes opening the same file and the relationship between the two different “task\_struct” structures is shown below.

Overall, we can say that “task\_struct” holds the data an operating system needs about a specific process. Among those data elements are: credentials, priority, PID (process ID), PPID (parent process ID), list of open resources, memory space range information, namespace information<sup>2</sup>, kprobes<sup>3</sup> instances and more.

Moreover, If you want to go over all of data elements I suggest going through the definition of “task\_struct” as part of the Linux source code<sup>4</sup>. Also, fun fact is that in kernel 6.2-rc1 “task\_struct” is referenced in 1398 files<sup>5</sup>.

Lastly, familiarity with “task\_struct” can help a lot with tracing and debugging tasks as shown in the online book “Dynamic Tracing with DTrace & SystemTap”<sup>6</sup>. Also, it is very handy when working with bpftrace. For example **sudo bpftrace -e 'kfunc:hrtimer\_wakeup { printf(“%s:%d\n”,curtask->comm,curtask->pid); }'**, which prints the pid and the process name of all processes calling the kernel function hrtimer\_wakeup<sup>7</sup>.



<sup>1</sup> <https://medium.com/@boutnaru/linux-processes-part-1-introduction-283f5b5b4197>

<sup>2</sup> <https://medium.com/system-weakness/linux-namespaces-part-1-dcee9c40fb68>

<sup>3</sup> <https://medium.com/@boutnaru/linux-instrumentation-part-2-kprobes-b089092c4cff>

<sup>4</sup> <https://elixir.bootlin.com/linux/v6.2-rc1/source/include/linux/sched.h#L737>

<sup>5</sup> [https://elixir.bootlin.com/linux/v6.2-rc1/A/ident/task\\_struct](https://elixir.bootlin.com/linux/v6.2-rc1/A/ident/task_struct)

<sup>6</sup> <https://myaut.github.io/dtrace-stap-book/kernel/proc.html>

<sup>7</sup> <https://medium.com/@boutnaru/the-linux-process-journey-pid-0-swapper-7868d1131316>

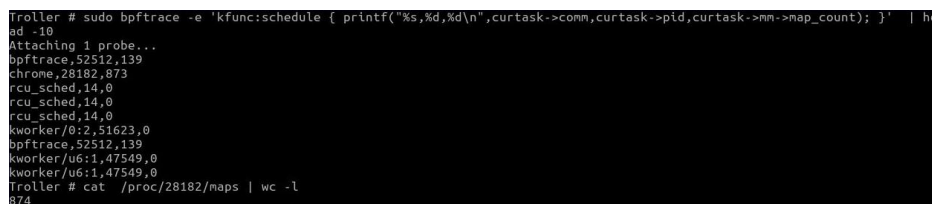
# struct mm\_struct

The goal of “mm\_struct” (aka “Memory Descriptor”) is to be used by the Linux kernel in order to represent the process’ address space<sup>8</sup>. It is the user-mode part which belongs to the task/process<sup>9</sup>. By the way, until kernel version 2.6.23 the struct was defined under “/include/linux/sched.h”<sup>10</sup>, since “2.6.24” it is located under “/include/linux/mm\_types.h”<sup>11</sup>. Also, there is a pointer from the process’ “task\_struct”<sup>12</sup> that refers to the address space of the process (“mm\_struct”) which is stored in the “mm” field<sup>13</sup>.

Overall, we can say that “mm\_struct” holds the data Linux needs about the memory address space of the process. Among those data elements are: “mm\_users” (the number of tasks using this address space), “map\_count” (the number of virtual memory areas, VMAs, used by the task) and “total\_vm” (the total number of pages mapped by the task). You can see part of the information stored in “mm\_struct” by going over “/proc/[PID]/maps” (“man proc”).

Moreover, If you want to go over all of data elements I suggest going through the definition of “mm\_struct” as part of the Linux source code<sup>14</sup>. Fun fact is that in kernel 6.2-rc1 “task\_struct” is referenced in 656 files<sup>15</sup>. Based on LXR<sup>16</sup> it seems that “mm\_struct” was added from kernel version 1.1.11<sup>17</sup> as we don’t see it in previous versions<sup>18</sup>.

Lastly, let us go over an example using bpftrace. We can use the following command: **sudo bpftrace -e 'kfunc:schedule { printf("%s,%d,%d\n",curtask->comm,curtask->pid,curtask->mm->map\_count); }'**. As shown in the screenshot below, the command prints the name, pid and VMA count for the current task every time the scheduler function is triggered. From the screenshot we can see that kernel threads have a count of zero VMAs (in their case current->mm==NULL). Also, we can see that the “map\_count” is one less<sup>19</sup> than the number of rows in “/proc/[PID]/maps”.



```
Trollier # sudo bpftrace -e 'kfunc:schedule { printf("%s,%d,%d\n",curtask->comm,curtask->pid,curtask->mm->map_count); }' | head -10
Attaching 1 probe...
bpftrace,52512,139
chrome,28182,873
rcu_sched,14,0
rcu_sched,14,0
rcu_sched,14,0
kworker/0:2,51623,0
bpftrace,52512,139
kworker/u6:1,47549,0
kworker/u6:1,47549,0
Trollier # cat /proc/28182/maps | wc -l
874
```

<sup>8</sup> <http://books.gigatux.nl/mirror/kerneldevelopment/0672327201/ch14lev1sec1.html>

<sup>9</sup> <https://medium.com/@boutnaru/linux-memory-management-part-1-introduction-896f376d3713>

<sup>10</sup> <https://elixir.bootlin.com/linux/v2.6.23/source/include/linux/sched.h#L369>

<sup>11</sup> [https://elixir.bootlin.com/linux/v2.6.24/source/include/linux/mm\\_types.h#L156](https://elixir.bootlin.com/linux/v2.6.24/source/include/linux/mm_types.h#L156)

<sup>12</sup> <https://medium.com/@boutnaru/linux-kernel-task-struct-829f51d97275>

<sup>13</sup> <https://elixir.bootlin.com/linux/v6.2-rc1/source/include/linux/sched.h#L870>

<sup>14</sup> [https://elixir.bootlin.com/linux/v6.2-rc1/source/include/linux/mm\\_types.h#L601](https://elixir.bootlin.com/linux/v6.2-rc1/source/include/linux/mm_types.h#L601)

<sup>15</sup> [https://elixir.bootlin.com/linux/v6.2-rc1/C/ident/mm\\_struct](https://elixir.bootlin.com/linux/v6.2-rc1/C/ident/mm_struct)

<sup>16</sup> <https://elixir.bootlin.com>

<sup>17</sup> <https://elixir.bootlin.com/linux/1.1.11/source/include/linux/sched.h#L214>

<sup>18</sup> [https://elixir.bootlin.com/linux/1.1.10/A/ident/mm\\_struct](https://elixir.bootlin.com/linux/1.1.10/A/ident/mm_struct)

<sup>19</sup> It is due to the mechanism of vsyscall

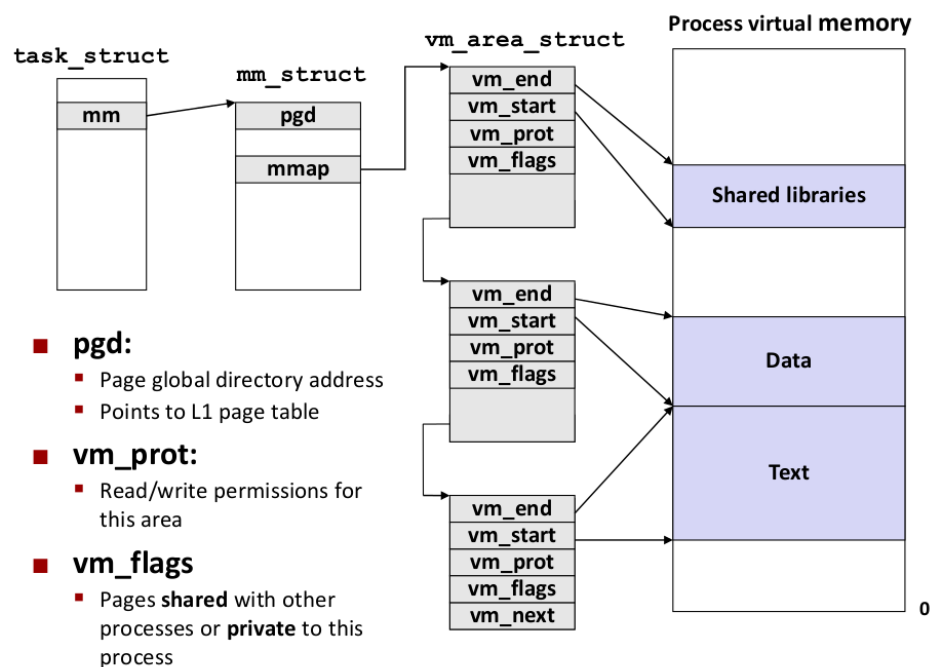
# struct vm\_area\_struct

“vm\_area\_struct” represents a contiguous memory area in a process's address space (virtual memory areas) - as shown in the diagram below<sup>20</sup>. It is used to track the permissions, properties, and operations associated with each memory area<sup>21</sup>.

This struct defines a memory VMM memory area. There is one of these per VM-area/task. A VM area is any part of the process virtual memory space that has a special rule for the page-fault handler. Think about shared libraries and executable area<sup>22</sup>.

Until kernel version 2.6.21 (including) “struct vm\_area\_struct” is defined in “/include/linux/mm.h”<sup>23</sup>. From kernel versions about it is defined in “/include/linux/mm\_types.h”<sup>24</sup>.

Lastly, by using “/proc/[PID]/maps” we can read the mapped regions and their access permissions (when using the mmap system call). For each region we can get the information about: its address range, pathname (in case mapped from a file), offset (in case mapped from a file), device (in case mapped from a file), inode (in case mapped from a file and permissions)<sup>25</sup>.



<sup>20</sup> [https://don7hao.github.io/2015/01/28/kernel/mm\\_struct/](https://don7hao.github.io/2015/01/28/kernel/mm_struct/)

<sup>21</sup> <http://books.gigatux.nl/mirror/kerneldevelopment/0672327201/ch14lev1sec2.html>

<sup>22</sup> <https://elixir.bootlin.com/linux/v2.6.21/source/include/linux/mm.h#L55>

<sup>23</sup> <https://elixir.bootlin.com/linux/v2.6.21/source/include/linux/mm.h#L60>

<sup>24</sup> [https://elixir.bootlin.com/linux/v6.5-rc1/source/include/linux/mm\\_types.h#L490](https://elixir.bootlin.com/linux/v6.5-rc1/source/include/linux/mm_types.h#L490)

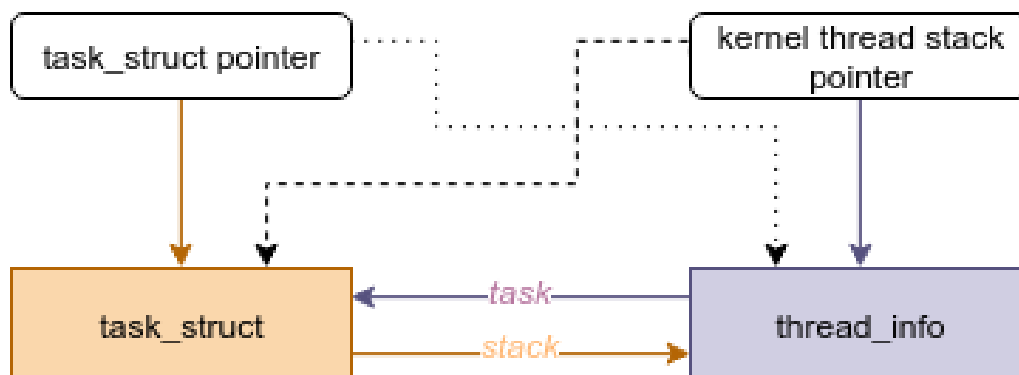
<sup>25</sup> <https://man7.org/linux/man-pages/man5/proc.5.html>

## struct thread\_info

“struct thread\_info” is a common low-level thread information accessors<sup>26</sup> (think about flags like signal pending, 32 address space on 64 bit, CPUID is not accessible in user mode and more). When “CONFIG\_THREAD\_INFO\_IN\_TASK” is defined, “struct thread\_info” is the first element of the “struct task\_struct”<sup>27</sup>. This means that each task has its own “struct thread\_info”.

Moreover, until kernel version 4.8 (including it) “struct thread\_info” contained a pointer to “struct task\_struct”<sup>28</sup> - the old relationship is shown in the diagram below<sup>29</sup>. But because it wasted too much space to keep it like that. By putting “struct thread\_info” in the start of “struct task\_struct” it makes getting from the “kernel stack”->“struct task\_struct” and from “struct task\_struct”->“kernel stack” very easy<sup>30</sup>. The first is done using the “current\_thread\_info” macro<sup>31</sup>. It uses “current\_task” which is a per-cpu variable<sup>32</sup>.

Lastly, it is important to understand that this data structure is CPU dependent and thus it is defined in the source code in the following location “/arch/[CPU\_ARCH]/include/asm/thread\_info.h”. Examples for “CPU\_ARCH” could be: “x86”<sup>33</sup>, “mips”<sup>34</sup>, “riscv”<sup>35</sup>, “arm64”<sup>36</sup> and more.



<sup>26</sup> [https://elixir.bootlin.com/linux/v6.5-rc3/source/include/linux/thread\\_info.h#L2](https://elixir.bootlin.com/linux/v6.5-rc3/source/include/linux/thread_info.h#L2)

<sup>27</sup> <https://medium.com/@boutnaru/linux-kernel-task-struct-829f51d97275>

<sup>28</sup> [https://elixir.bootlin.com/linux/v4.8.17/source/arch/x86/include/asm/thread\\_info.h#L56](https://elixir.bootlin.com/linux/v4.8.17/source/arch/x86/include/asm/thread_info.h#L56)

<sup>29</sup> <https://blog.lexfo.fr/cve-2017-11176-linux-kernel-exploitation-part4.html>

<sup>30</sup> <https://stackoverflow.com/questions/61886139/why-thread-info-should-be-the-first-element-in-task-struct>

<sup>31</sup> [https://elixir.bootlin.com/linux/v6.5-rc3/source/include/linux/thread\\_info.h#L24](https://elixir.bootlin.com/linux/v6.5-rc3/source/include/linux/thread_info.h#L24)

<sup>32</sup> <https://elixir.bootlin.com/linux/v6.5-rc3/source/arch/x86/include/asm/current.h#L41>

<sup>33</sup> [https://elixir.bootlin.com/linux/v6.5-rc3/source/arch/x86/include/asm/thread\\_info.h#L56](https://elixir.bootlin.com/linux/v6.5-rc3/source/arch/x86/include/asm/thread_info.h#L56)

<sup>34</sup> [https://elixir.bootlin.com/linux/v6.5-rc3/source/arch/mips/include/asm/thread\\_info.h#L25](https://elixir.bootlin.com/linux/v6.5-rc3/source/arch/mips/include/asm/thread_info.h#L25)

<sup>35</sup> [https://elixir.bootlin.com/linux/v6.5-rc3/source/arch/riscv/include/asm/thread\\_info.h#L51](https://elixir.bootlin.com/linux/v6.5-rc3/source/arch/riscv/include/asm/thread_info.h#L51)

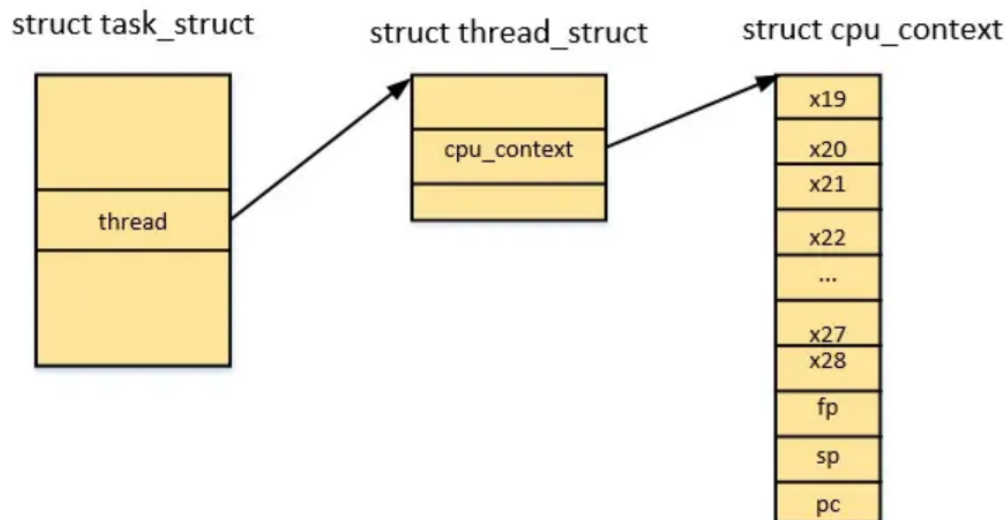
<sup>36</sup> [https://elixir.bootlin.com/linux/v6.5-rc3/source/arch/arm64/include/asm/thread\\_info.h#L24](https://elixir.bootlin.com/linux/v6.5-rc3/source/arch/arm64/include/asm/thread_info.h#L24)

## struct thread\_struct

Overall, the goal of “struct thread\_struct” which hold CPU-specific state of a task<sup>37</sup>. Among the information “struct thread\_struct” holds we can include things like page fault information (like the address that caused the page fault and fault code). Also, it can include a set of registers of the current CPU - as shown in the the diagram below<sup>38</sup>.

On x86 the variable which is part of “struct task\_struct”<sup>39</sup> must be at the end of the struct. The reason for that is it contain a variable-sized structure<sup>40</sup>.

Moreover, like “struct thread\_info”<sup>41</sup> also “struct thread\_struct” is CPU/Architecture dependent and thus it is defined in the source code in the following location “arch/[CPU\_ARCH]/include/asm/processor.h”<sup>42</sup>. For example x86<sup>43</sup> and arm64<sup>44</sup>.



<sup>37</sup> <https://elixir.bootlin.com/linux/v6.5-rc4/source/include/linux/sched.h#L1540>

<sup>38</sup> <https://kernel.0voice.com/forum.php?mod=viewthread&tid=2920>

<sup>39</sup> <https://medium.com/@boutnaru/linux-kernel-task-struct-829f51d97275>

<sup>40</sup> <https://elixir.bootlin.com/linux/v6.5-rc4/source/include/linux/sched.h#L1544>

<sup>41</sup> <https://medium.com/@boutnaru/the-linux-kernel-data-structure-journey-struct-thread-info-4e70bc20d279>

<sup>42</sup> [https://elixir.bootlin.com/linux/v6.5-rc4/C/ident/thread\\_struct](https://elixir.bootlin.com/linux/v6.5-rc4/C/ident/thread_struct)

<sup>43</sup> <https://elixir.bootlin.com/linux/v6.5-rc4/source/arch/x86/include/asm/processor.h#L414>

<sup>44</sup> <https://elixir.bootlin.com/linux/v6.5-rc4/source/arch/arm64/include/asm/processor.h#L147>

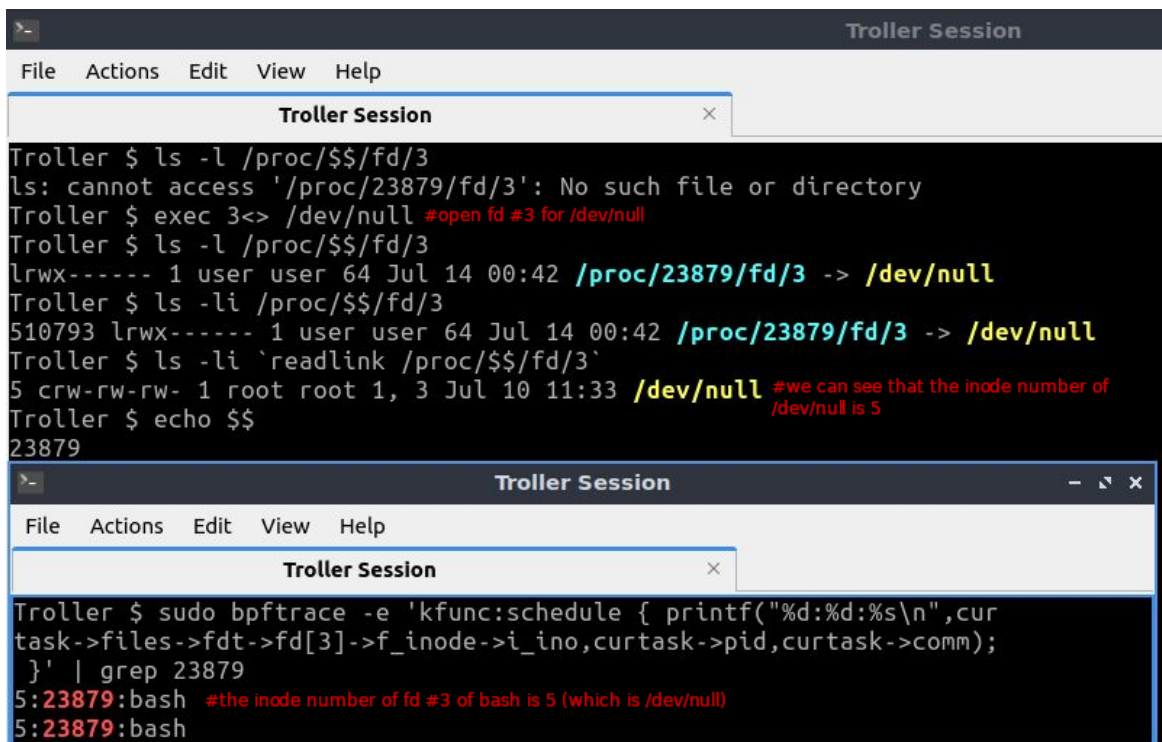


# struct inode

Overall, the goal of an inode<sup>45</sup> object is to hold all the data needed by the kernel in order to perform actions on directories/files. In case of Unix/Linux style filesystem (like ext4) it can read from the on-disk inode<sup>46</sup>. Moreover, inode is an essential component of Unix filesystem and of VFS (Virtual File system). We can say that inode exists both in memory as part of VFS and as an entity on disk. The inode as part of VFS is represented by “struct inode”<sup>47</sup>.

Also, the “struct inode” is defined in “/include/linux/fs.h” in the source tree of Linux<sup>48</sup>. It contains different members such as: jiffies of first dirtying, LRU list, backing dev writeback list, inode number, a mask of notify events that the inode cares about and more.

Lastly, let us go over an example using bpftrace - as shown in the screenshot below. First we open fd 3 to “/dev/null” (using exec). Then we check and see that the inode of “/dev/null” is 5. We can use the following command: **sudo bpftrace -e 'kfunc:schedule { printf("%d:%d:%s\n",curtask->files->fdt->fd[3]->f\_inode->i\_ino,curtask->pid,curtask->comm); }'**. As shown in the screenshot below it goes over different structs until getting to the inode of fd(3) and prints its number (which is 5). By the way, the example was conducted on kernel version 5.15.



```
Troller $ ls -l /proc/$$/fd/3
ls: cannot access '/proc/23879/fd/3': No such file or directory
Troller $ exec 3<> /dev/null #open fd #3 for /dev/null
Troller $ ls -l /proc/$$/fd/3
lrwx----- 1 user user 64 Jul 14 00:42 /proc/23879/fd/3 -> /dev/null
Troller $ ls -li /proc/$$/fd/3
510793 lrwx----- 1 user user 64 Jul 14 00:42 /proc/23879/fd/3 -> /dev/null
Troller $ ls -li `readlink /proc/$$/fd/3`
5 crw-rw-rw- 1 root root 1, 3 Jul 10 11:33 /dev/null #we can see that the inode number of /dev/null is 5
Troller $ echo $$
23879
```

```
Troller $ sudo bpftrace -e 'kfunc:schedule { printf("%d:%d:%s\n",curtask->files->fdt->fd[3]->f_inode->i_ino,curtask->pid,curtask->comm); }' | grep 23879
5:23879:bash #the inode number of fd #3 of bash is 5 (which is /dev/null)
5:23879:bash
```

<sup>45</sup> <https://medium.com/@boutnaru/linux-what-is-an-inode-7ba47a519940>

<sup>46</sup> <http://books.gigatux.nl/mirror/kerneldevelopment/0672327201/ch12lev1sec6.html>

<sup>47</sup> [https://linux-kernel-labs.github.io/refs/heads/master/labs/filesystems\\_part2.html](https://linux-kernel-labs.github.io/refs/heads/master/labs/filesystems_part2.html)

<sup>48</sup> <https://elixir.bootlin.com/linux/v6.5-rc1/source/include/linux/fs.h#L608>