

[about](#) [summary](#) [refs](#) [log](#) [tree](#) [commit](#) [diff](#) [stats](#)[log msg](#) 

author Darrick J. Wong <djwong@kernel.org> 2023-02-16 10:55:48 -0800  
committer Greg Kroah-Hartman <gregkh@linuxfoundation.org> 2023-03-17 08:30:37 +0100  
commit [a70b49dc7eee5dbe3775a650ce598e3557ff5475](#) ([patch](#))  
tree [e77dba1394f748bcc1f667bd1da605533806e04d](#)  
parent [e49ddab12b0a19ccbcfc223c1d3976dcb253b063](#) ([diff](#))  
download [linux-a70b49dc7eee5dbe3775a650ce598e3557ff5475.tar.gz](#)

**diff options**

context:    
space:    
mode:

**ext4: fix another off-by-one fsmap error on 1k block filesystems**

commit [c993799baf9c5861f8df91beb80e1611b12efcbd](#) upstream.

Apparently syzbot figured out that issuing this FSMAP call:

```
struct fsmap_head cmd = {
    .fmh_count      = ...;
    .fmh_keys       = {
        { .fmr_device = /* ext4 dev */, .fmr_physical = 0, },
        { .fmr_device = /* ext4 dev */, .fmr_physical = 0, },
    },
...
};
ret = ioctl(fd, FS_IOC_GETFSMAP, &cmd);
```

Produces this crash if the underlying filesystem is a 1k-block ext4 filesystem:

```
kernel BUG at fs/ext4/ext4.h:3331!
invalid opcode: 0000 [#1] PREEMPT SMP
CPU: 3 PID: 3227965 Comm: xfs_io Tainted: G          W 0      6.2.0-rc8-achx
Hardware name: QEMU Standard PC (Q35 + ICH9, 2009), BIOS 1.15.0-1 04/01/2014
RIP: 0010:ext4_mb_load_buddy_gfp+0x47c/0x570 [ext4]
RSP: 0018:fffffc90007c03998 EFLAGS: 00010246
RAX: fffff888004978000 RBX: fffffc90007c03a20 RCX: fffff888041618000
RDX: 0000000000000000 RSI: 0000000000005a4 RDI: ffffffff0c99b11
RBP: fffff888012330000 R08: ffffffa0c2b7d0 R09: 0000000000000400
R10: fffffc90007c03950 R11: 0000000000000000 R12: 0000000000000001
R13: 00000000ffffffffff R14: 0000000000000c40 R15: fffff88802678c398
FS: 00007fdf2020c880(0000) GS:fffff88807e100000(0000) knlGS:0000000000000000
CS: 0010 DS: 0000 ES: 0000 CR0: 0000000080050033
CR2: 00007ffd318a5fe8 CR3: 000000007f80f001 CR4: 00000000001706e0
Call Trace:
<TASK>
ext4_mballloc_query_range+0x4b/0x210 [ext4 dfa189daddffe8fec3cdfd00564e0f265a8ab80]
ext4_getfsmap_datadev+0x713/0x890 [ext4 dfa189daddffe8fec3cdfd00564e0f265a8ab80]
ext4_getfsmap+0x2b7/0x330 [ext4 dfa189daddffe8fec3cdfd00564e0f265a8ab80]
ext4_ioc_getfsmap+0x153/0x2b0 [ext4 dfa189daddffe8fec3cdfd00564e0f265a8ab80]
__ext4_ioctl+0x2a7/0x17e0 [ext4 dfa189daddffe8fec3cdfd00564e0f265a8ab80]
__x64_sys_ioctl+0x82/0xa0
do_syscall_64+0x2b/0x80
entry_SYSCALL_64_after_hwframe+0x46/0xb0
```

```
RIP: 0033:0x7fdf20558aff
RSP: 002b:00007ffd318a9e30 EFLAGS: 00000246 ORIG_RAX: 0000000000000010
RAX: ffffffff000000000000000000000000 RBX: 000000000000200c0 RCX: 00007fdf20558aff
RDX: 00007fdf1feb2010 RSI: 00000000c0c0583b RDI: 000000000000000000000003
RBP: 00005625c0634be0 R08: 00005625c0634c40 R09: 000000000000000000000001
R10: 0000000000000000 R11: 00000000000000246 R12: 00007fdf1feb2010
R13: 00005625be70d994 R14: 00000000000000800 R15: 000000000000000000000000
```

For GETFSMAP calls, the caller selects a physical block device by writing its block number into fsmap\_head.fmh\_keys[0].fmr\_device. To query mappings for a subrange of the device, the starting byte of the range is written to fsmap\_head.fmh\_keys[0].fmr\_physical and the last byte of the range goes in fsmap\_head.fmh\_keys[1].fmr\_physical.

IOWs, to query what mappings overlap with bytes 3-14 of /dev/sda, you'd set the inputs as follows:

```
fmh_keys[0] = { .fmr_device = major(8, 0), .fmr_physical = 3},
fmh_keys[1] = { .fmr_device = major(8, 0), .fmr_physical = 14},
```

Which would return you whatever is mapped in the 12 bytes starting at physical offset 3.

The crash is due to insufficient range validation of keys[1] in ext4\_getfsmap\_datadev. On 1k-block filesystems, block 0 is not part of the filesystem, which means that s\_first\_data\_block is nonzero. ext4\_get\_group\_no\_and\_offset subtracts this quantity from the blocknr argument before cracking it into a group number and a block number within a group. IOWs, block group 0 spans blocks 1-8192 (1-based) instead of 0-8191 (0-based) like what happens with larger blocksizes.

The net result of this encoding is that blocknr < s\_first\_data\_block is not a valid input to this function. The end\_fsb variable is set from the keys that are copied from userspace, which means that in the above example, its value is zero. That leads to an underflow here:

```
blocknr = blocknr - le32_to_cpu(es->s_first_data_block);
```

The division then operates on -1:

```
offset = do_div(blocknr, EXT4_BLOCKS_PER_GROUP(sb)) >>
EXT4_SB(sb)->s_cluster_bits;
```

Leaving an impossibly large group number ( $2^{32}-1$ ) in blocknr. ext4\_getfsmap\_check\_keys checked that keys[0].fmr\_physical and keys[1].fmr\_physical are in increasing order, but ext4\_getfsmap\_datadev adjusts keys[0].fmr\_physical to be at least s\_first\_data\_block. This implies that we have to check it again after the adjustment, which is the piece that I forgot.

Reported-by: syzbot+6be2b977c89f79b6b153@syzkaller.appspotmail.com  
Fixes: 4a4956249dac ("ext4: fix off-by-one fsmap error on 1k block filesystems")  
Link: <https://syzkaller.appspot.com/bug?id=79d5768e9bfe362911ac1a5057a36fc6b5c30002>  
Cc: stable@vger.kernel.org  
Signed-off-by: Darrick J. Wong <djhong@kernel.org>  
Link: <https://lore.kernel.org/r/Y+58NPTH7VNGgzdd@magnolia>  
Signed-off-by: Theodore Ts'o <tytso@mit.edu>  
Signed-off-by: Greg Kroah-Hartman <gregkh@linuxfoundation.org>

## Diffstat

-rw-r--r--	fs/ext4/fsmmap.c	2
------------	------------------	---

1 files changed, 2 insertions, 0 deletions

```
diff --git a/fs/ext4/fsmap.c b/fs/ext4/fsmap.c
index 1a4d42a1b161de..f1ad06e688efc4 100644
--- a/fs/ext4/fsmap.c
+++ b/fs/ext4/fsmap.c
@@ -499,6 +499,8 @@ static int ext4_getfsmap_datadev(struct super_block *sb,
        keys[0].fmr_physical = bofs;
        if (keys[1].fmr_physical >= eofsl)
            keys[1].fmr_physical = eofsl - 1;
+       if (keys[1].fmr_physical < keys[0].fmr_physical)
+
            return 0;
        start_fsb = keys[0].fmr_physical;
        end_fsb = keys[1].fmr_physical;
```

---

generated by cgit 1.2.3-korg (git 2.43.0) at 2025-05-03 16:43:24 +0000