



author Omar Sandoval <osandov@fb.com> 2025-04-25 01:51:24 -0700  
committer Greg Kroah-Hartman <gregkh@linuxfoundation.org> 2025-05-02 08:01:44 +0200  
commit 50a665496881262519f115f1bfe5822f30580eb0 (patch)  
tree bfe74104ab3b663f3be4b5a0483b502a807e5154  
parent ea1fd673d95eca322d2f4fc1de6d41280a13a0ee (diff)  
download linux-50a665496881262519f115f1bfe5822f30580eb0.tar.gz

**diff options**

context:  space:  mode:

**sched/eevdf: Fix se->slice being set to U64\_MAX and resulting crash**

[ Upstream commit bbce3de72be56e4b5f68924b7da9630cc89aa1a8 ]

There is a code path in dequeue\_entities() that can set the slice of a sched\_entity to U64\_MAX, which sometimes results in a crash.

The offending case is when dequeue\_entities() is called to dequeue a delayed group entity, and then the entity's parent's dequeue is delayed. In that case:

1. In the if (entity\_is\_task(se)) else block at the beginning of dequeue\_entities(), slice is set to cfs\_rq\_min\_slice(group\_cfs\_rq(se)). If the entity was delayed, then it has no queued tasks, so cfs\_rq\_min\_slice() returns U64\_MAX.
2. The first for\_each\_sched\_entity() loop dequeues the entity.
3. If the entity was its parent's only child, then the next iteration tries to dequeue the parent.
4. If the parent's dequeue needs to be delayed, then it breaks from the first for\_each\_sched\_entity() loop \_without updating slice\_.
5. The second for\_each\_sched\_entity() loop sets the parent's ->slice to the saved slice, which is still U64\_MAX.

This throws off subsequent calculations with potentially catastrophic results. A manifestation we saw in production was:

6. In update\_entity\_lag(), se->slice is used to calculate limit, which ends up as a huge negative number.
7. limit is used in se->vlag = clamp(vlag, -limit, limit). Because limit is negative, vlag > limit, so se->vlag is set to the same huge negative number.
8. In place\_entity(), se->vlag is scaled, which overflows and results in another huge (positive or negative) number.
9. The adjusted lag is subtracted from se->vruntime, which increases or decreases se->vruntime by a huge number.
10. pick\_eevdf() calls entity\_eligible()/vruntime\_eligible(), which incorrectly returns false because the vruntime is so far from the other vruntimes on the queue, causing the (vruntime - cfs\_rq->min\_vruntim) \* load calculation to overflow.
11. Nothing appears to be eligible, so pick\_eevdf() returns NULL.
12. pick\_next\_entity() tries to dereference the return value of pick\_eevdf() and crashes.

Dumping the cfs\_rq states from the core dumps with drgn showed tell-tale huge vruntime ranges and bogus vlag values, and I also traced se->slice being set to U64\_MAX on live systems (which was usually "benign" since

the rest of the runqueue needed to be in a particular state to crash).

Fix it in `dequeue_entities()` by always setting `slice` from the first non-empty `cfs_rq`.

Fixes: aef6987d8954 ("sched/eevdf: Propagate `min_slice` up the cgroup hierarchy")

Signed-off-by: Omar Sandoval <osandov@fb.com>

Signed-off-by: Peter Zijlstra (Intel) <peterz@infradead.org>

Signed-off-by: Ingo Molnar <mingo@kernel.org>

Link: <https://lkml.kernel.org/r/f0c2d1072be229e1bddc73c0703919a8b00c652.1745570998.git.osandov@fb.com>

Signed-off-by: Sasha Levin <sashal@kernel.org>

## Diffstat

-rw-r--r-- kernel/sched/fair.c 4
----------------------------------

1 files changed, 1 insertions, 3 deletions

**diff --git a/kernel/sched/fair.c b/kernel/sched/fair.c**

**index 89c7260103e18b..3d9b68a347b764 100644**

**---** a/**kernel/sched/fair.c**

**+++ b/**kernel/sched/fair.c****

**@@ -7083,9 +7083,6 @@ static int dequeue\_entities(struct rq \*rq, struct sched\_entity \*se, int flags)**

```
        h_nr_idle = task_has_idle_policy(p);
        if (task_sleep || task_delayed || !se->sched_delayed)
            h_nr_runnable = 1;
```

**- } else {**

```
-        cfs_rq = group_cfs_rq(se);
-        slice = cfs_rq_min_slice(cfs_rq);
-    }
```

```
    for_each_sched_entity(se) {
```

**@@ -7095,6 +7092,7 @@ static int dequeue\_entities(struct rq \*rq, struct sched\_entity \*se, int flags)**

```
        if (p && &p->se == se)
            return -1;
```

```
+        slice = cfs_rq_min_slice(cfs_rq);
        break;
}
```