

---

Subject: Re: [BUG]: Crash with CONFIG\_FAIR\_CGROUP\_SCHED=y  
Posted by [Sukadev Bhattiprolu](#) on Sat, 10 Nov 2007 23:13:47 GMT  
[View Forum Message](#) <> [Reply to Message](#)

---

Serge E. Hallyn [serue@us.ibm.com] wrote:

| Quoting Srivatsa Vaddagiri (vatsa@linux.vnet.ibm.com):  
| > On Fri, Nov 09, 2007 at 09:45:21AM +0100, Dmitry Adamushko wrote:  
| > > Humm... the 'current' is not kept within the tree but  
| > > current->se.on\_rq is supposed to be '1',  
| > > so the old code looks ok to me (at least for the 'leaf' elements).  
| >  
| > You are damned right! Sorry my mistake with the previous analysis and  
| > (as I now find out) testing :(  
| >  
| > There are couple of problems discovered by Suka's test:  
| >  
| > - The test requires the cgroup filesystem to be mounted with  
| > at least the cpu and ns options (i.e both namespace and cpu  
| > controllers are active in the same hierarchy).  
| >  
| > # mkdir /dev/cpuctl  
| > # mount -t cgroup -ocpu,ns none cpuctl  
| > (or simply)  
| > # mount -t cgroup none cpuctl -> Will activate all controllers  
| > in same hierarchy.  
| >  
| > - The test invokes clone() with CLONE\_NEWNS set. This causes a a new child  
| > to be created, also a new group (do\_fork->copy\_namespaces->ns\_cgroup\_clone->  
| > cgroup\_clone) and the child is attached to the new group (cgroup\_clone->  
| > attach\_task->sched\_move\_task). At this point in time, the child's scheduler  
| > related fields are uninitialized (including its on\_rq field, which it has  
| > inherited from parent). As a result sched\_move\_task thinks its on  
| > runqueue, when it isn't.  
| >  
| > As a solution to this problem, I moved sched\_fork() call, which  
| > initializes scheduler related fields on a new task, before  
| > copy\_namespaces(). I am not sure though whether moving up will  
| > cause other side-effects. Do you see any issue?  
| >  
| > - The second problem exposed by this test is that task\_new\_fair()  
| > assumes that parent and child will be part of the same group (which  
| > needn't be as this test shows). As a result, cfs\_rq->curr can be NULL  
| > for the child.  
| >  
| > The solution is to test for curr pointer being NULL in  
| > task\_new\_fair().  
| >  
| >

| > With the patch below, I could run ns\_exec() fine w/o a crash.  
| >  
| > Suka, can you verify whether this patch fixes your problem?  
|  
| Works on my machine. Thanks!

And mine too. Thanks,

|  
| > --  
| >  
| > Fix copy\_namespace() <-> sched\_fork() dependency in do\_fork, by moving  
| > up sched\_fork().  
| >  
| > Also introduce a NULL pointer check for 'curr' in task\_new\_fair().  
| >  
| > Signed-off-by : Srivatsa Vaddagiri <vatsa@linux.vnet.ibm.com>  
|  
| Tested-by: Serge Hallyn <serue@us.ibm.com>  
Tested-by: Sukadev Bhattiprolu <sukadev@us.ibm.com>

---

Containers mailing list  
Containers@lists.linux-foundation.org  
<https://lists.linux-foundation.org/mailman/listinfo/containers>

---