

---

Subject: Re: [PATCH 14/15] Destroy pid namespace on init's death

Posted by [Pavel Emelianov](#) on Mon, 30 Jul 2007 11:56:50 GMT

[View Forum Message](#) <> [Reply to Message](#)

---

Oleg Nesterov wrote:

> On 07/26, Pavel Emelyanov wrote:

>> @@ -895,6 +915,7 @@ fastcall NORET\_TYPE void do\_exit(long co

>> {

>> struct task\_struct \*tsk = current;

>> int group\_dead;

>> + struct pid\_namespace \*pid\_ns = tsk->nsproxy->pid\_ns;

>>

>> profile\_task\_exit(tsk);

>>

>> @@ -905,9 +926,10 @@ fastcall NORET\_TYPE void do\_exit(long co

>> if (unlikely(!tsk->pid))

>> panic("Attempted to kill the idle task!");

>> if (unlikely(tsk == task\_child\_reaper(tsk))) {

>> - if (task\_active\_pid\_ns(tsk) != &init\_pid\_ns)

>> - task\_active\_pid\_ns(tsk)->child\_reaper =

>> - init\_pid\_ns.child\_reaper;

>> + if (pid\_ns != &init\_pid\_ns) {

>> + zap\_pid\_ns\_processes(pid\_ns);

>> + pid\_ns->child\_reaper = init\_pid\_ns.child\_reaper;

>> + }

>> else

>> panic("Attempted to kill init!");

>

> No, no, this is wrong. Yes, the current code is buggy too, I'll send

> the fix.

>

> I think this code should be moved below under the "if (group\_dead)",

> and we should use tsk->group\_leader.

>

>> +void zap\_pid\_ns\_processes(struct pid\_namespace \*pid\_ns)

>> +{

>> + int i;

>> + int nr;

>> + int nfree;

>> + int options = WNOHANG|WEXITED|\_\_WALL;

>> +

>> +repeat:

>> + /\*

>> + \* We know pid == 1 is terminating. Find remaining pid\_ts

>> + \* in the namespace, signal them and then wait for them

>> + \* exit.

>> + \*/

>> + nr = next\_pidmap(pid\_ns, 1);

```

>> + while (nr > 0) {
>> +   kill_proc_info(SIGKILL, SEND_SIG_PRIV, nr);
>> +   nr = next_pidmap(pid_ns, nr);
>> + }
>> +
>> + nr = next_pidmap(pid_ns, 1);
>> + while (nr > 0) {
>> +   do_wait(nr, options, NULL, NULL, NULL);
>
> When the first child of init exits, it sends SIGCHLD. After that,
> do_wait() will never sleep, so we are doing a busy-wait loop.
> Not good, especially when we have a niced child, can livelock.
>
>> + nr = next_pidmap(pid_ns, nr);
>> + }
>> +
>> + nfree = 0;
>> + for (i = 0; i < PIDMAP_ENTRIES; i++)
>> +   nfree += atomic_read(&pid_ns->pidmap[i].nr_free);
>> +
>> + /*
>> +  * If pidmap has entries for processes other than 0 and 1, retry.
>> +  */
>> + if (nfree < (BITS_PER_PAGE * PIDMAP_ENTRIES - 2))
>> +   goto repeat;
>
> This doesn't look right.
>
> Suppose that some "struct pid" was pinned from the parent namespace.
> In that case zap_pid_ns_processes() will burn CPU until put_pid(), bad.

```

Nope. struct pid can be pinned, but the pidmap "fingerprint" cannot.  
 So as soon as the release\_task() is called the pidmap becomes free and  
 we can proceed.

However I agree with the "burn CPU" issue - wait must sleep if needed.

```

> I think we can rely on forget_original_child() and do something like
> this:
>
> zap_active_ns_processes(void)
> {
>   // kill all tasks in our ns and below
>   kill(-1, SIGKILL);

```

That would be too slow to walk through all the tasks in a node searching  
 for a couple of them we need. find\_get\_pid() looks better to me.

```
> do {  
>   clear_thread_flag(TIF_SIGPENDING);  
> } while (wait(NULL) != -ECHLD);  
> }  
>  
> Oleg.  
>  
>
```

Thanks,  
Pavel

---