

Pavel Emelianov <xemul@openvz.org> writes:

> Eric W. Biederman wrote:
>> Pavel Emelianov <xemul@openvz.org> writes:
>>
>>> That's how OpenVZ sees the pid namespaces.
>>>
>>> The main idea is that kernel keeps operating with tasks pid
>>> as it did before, but each task obtains one more pid for each
>>> pid type - the virtual pid. When putting the pid to user or
>>> getting the pid from it kernel operates with the virtual ones.
>>
>> Just a quick reaction.
>>
>> - I would very much like to see a minimum of 3 levels of pids,
>
> Why not 4? From my part, I would like to know, why such nesting
> is important. We have plain IPC namespaces and nobody cares.
> We will have isolated network namespaces, why pids are exception?

4+ is fine, and something we will probably care about someday.
3 seems to be the minimum necessary to get people thinking about
adding more so we don't have arbitrary special cases, especially
in the user interface. At 3 the things are simple enough we don't
have to allocate additional data structures etc.

If we don't need nesting we don't even need 2 levels, and we
can remove the global pid. But we have had that conversation
and especially for the current OpenVZ usage we need nesting.

Having more then two layers means we are prepared to use pid namespaces more
generally. It really isn't that much harder.

>> being supported. Otherwise it is easy to overlook some of the
>> cases that are required to properly support nesting, which long
>> terms seems important.
>>
>> - Semantically fork is easier then unshare. Unshare can mean
>
> This is not. When you fork, the kid shares the session and the
> group with its parent, but moving this pids to new ns is bad - the
> parent will happen to be half-moved. Thus you need to break the
> session and the group in fork(), but this is extra complexity.

Nope. You will just need to have the child call `setsid()` if you don't want to share the session and the group.

You can perfectly well share the sid and group with the parent, because internal to the kernel pids aren't numeric, they are struct pid pointers.

There is the question of do you use foreign pid handling to display the session and the group, or do you allocate pids for the session and the group in the new pid namespace. At this point foreign pid handling looks sufficient.

>> a lot of things, and it is easy to pick a meaning that has weird
>> side effects. Your implementation has a serious problem in that you
>> change the value of `getpid()` at runtime. Glibc does not know how to
>> cope with the value of `getpid()` changing.
>
> This pid changing happens only once per task lifetime.

Unshare isn't once per task lifetime, unless you added some extra constraints.

> Though I haven't
> seen any problems with glibc for many years running OpenVZ and I think,
> that if glibc will want to cache this `getpid()` value we can teach it to
> uncache this value in case someone called `unshare()` with `CLONE_NEWPIDS`.

glibc very much caches the results of `getpid()`.

If you want to teach glibc not to cache `getpid()` for free. The only way I know to get glibc to invalidate its pid cache is to call `fork`.

Eric
