

---

Subject: Re: [PATCH 1/9] Add a user\_namespace as creator/owner of uts\_namespace

Posted by [ebiederm](#) on Wed, 23 Feb 2011 21:21:28 GMT

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

---

David Howells <dhowells@redhat.com> writes:

> Serge E. Hallyn <serge@hallyn.com> wrote:

>  
>> struct uts\_namespace {  
>> struct kref kref;  
>> struct new\_utsname name;  
>> + struct user\_namespace \*user\_ns;  
>> };  
>  
> If a uts\_namespace belongs to a user\_namespace, should CLONE\_NEWUSER then  
> imply CLONE\_NEWUTS?  
>  
> Or is uts\_namespace::user\_ns more an implication that the set of users in  
> user\_namespace are the only ones authorised to alter the uts data.

The later.

> I presume that the uts\_namespace of a process must be owned by one of the  
> user\_namespaces in the alternating inheritance chain of namespaces and their  
> creators leading from current\_user\_ns() to init\_user\_ns.  
>  
> With that in mind, looking at patch 3:  
>  
> - if (!capable(CAP\_SYS\_ADMIN))  
> + if (!ns\_capable(current->nsproxy->uts\_ns->user\_ns, CAP\_SYS\_ADMIN))  
>  
> what is it you're actually asking? I presume it's 'does this user have  
> CAP\_SYS\_ADMIN capability over objects belonging to the uts\_namespace's  
> user\_namespace?'

Yes.

> So, to look at the important bit of patch 2:  
>  
> -int cap\_capable(struct task\_struct \*tsk, const struct cred \*cred, int cap,  
> - int audit)  
> +int cap\_capable(struct task\_struct \*tsk, const struct cred \*cred,  
> + struct user\_namespace \*targ\_ns, int cap, int audit)  
> {  
> - return cap\_raised(cred->cap\_effective, cap) ? 0 : -EPERM;  
> + for (;;) {  
> + /\* The creator of the user namespace has all caps. \*/

```

> + if (targ_ns != &init_user_ns && targ_ns->creator == cred->user)
> + return 0;
> +
> + /* Do we have the necessary capabilities? */
> + if (targ_ns == cred->user->user_ns)
> + return cap_raised(cred->cap_effective, cap) ? 0 : -EPERM;
> +
> + /* Have we tried all of the parent namespaces? */
> + if (targ_ns == &init_user_ns)
> + return -EPERM;
> +
> + /* If you have the capability in a parent user ns you have it
> + * in the over all children user namespaces as well, so see
> + * if this process has the capability in the parent user
> + * namespace.
> + */
> + targ_ns = targ_ns->creator->user_ns;
> + }
> +
> + /* We never get here */
> + return -EPERM;
> }
>
> On entry, as we're called from ns_capable(), cred->user is the user that the
> current process is running as, and, as such, may be in a separate namespace
> from uts_namespace - which may itself be in a separate namespace from
> init_user_ns.
>
> So, assume for the sake of argument that there are three user_namespaces along
> the chain from the calling process to the root, and that the uts_namespace
> belongs to the middle one.

```

So we have the nested stack of:

```

user_ns3->creator->user_ns == user_ns2
user_ns2->creator->user_ns == &init_user_ns
uts_ns2->user_ns == user_ns2

```

```

>
> if (targ_ns != &init_user_ns && targ_ns->creator == cred->user)
> return 0;
>
> Can never match because targ_ns->creator cannot be cred->user; even if the
> uts_namespace belongs to our namespace, given that the creator lies outside
> our namespace.

```

Initially we come in with `targ_ns == user_ns2` and `cred->user->user_ns` in one of (`user_ns3`, `user_ns2`, or `&init_user_ns`).

targ\_ns takes on values user\_ns2 and &init\_user\_ns.

So when targ\_ns becomes &init\_user\_ns. If the user in question is uts\_ns2->user\_ns->creator. This check will indeed match.

```
> if (targ_ns == cred->user->user_ns)
>   return cap_raised(cred->cap_effective, cap) ? 0 : -EPERM;
>
> Can only match if we are in the target user_namespace (ie. the one to which
> uts_namespace belongs), whether or not we have CAP_SYS_ADMIN.
```

As before targ\_ns takes on values of user\_ns2 and &init\_user\_ns.

Which means this check will match if we have CAP\_SYS\_ADMIN in &init\_user\_ns or in user\_ns2.

```
> Which means that unless the uts_namespace belongs to our user_namespace, we
> cannot change it. Is that correct?
```

No. If you are root in a parent namespace you can also change it.

```
> So ns_capable() restricts you to only doing interesting things to objects that
> belong to a user_namespace if they are in your own user_namespace. Is that
> correct?
```

No. Root outside your user namespace is also allowed to do interesting things.

```
> If that is so, is the loop required for ns_capable()?
```

Yes.

```
> Looking further at patch 2:
```

```
>
> #define nsown_capable(cap) (ns_capable(current_user_ns(), (cap)))
>
> Given what I've said above, I presume the loop isn't necessary here either.
>
>
> I think you're using ns_capable() in two different ways:
>
> (1) You're using it to see if a process has power over its descendents in a
>     user_namespace that can be traced back to a clone() that it did with
>     CLONE_NEWUSER.
>
>     For example, automatically granting a process permission to kill
>     descendents in a namespace it created.
```

>  
> (2) You're using it to see if a process can access objects that might be  
> outside its own user\_namespace.  
>  
> For example, setting the hostname.  
>  
> Is it worth giving two different interfaces to make this clearer (even if they  
> actually do the same thing)?  
>  
>  
> Sorry if this seems rambling, but I'm trying to get my head round your  
> code.

I am all for making that loop a little clearer, because it is something you have to pause and think about to understand but so far I don't think the loop is wrong, and it is simple.

Eric

---

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

---