
Subject: Re: [RFC][PATCH 1/2] add user namespace [try #2]

Posted by [dev](#) on Thu, 07 Sep 2006 15:59:24 GMT

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

comments below

> This patch adds the user namespace.
>
> Basically, it allows a process to unshare its user_struct table,
> resetting at the same time its own user_struct and all the associated
> accounting.
>
> A new root user (uid == 0) is added to the user namespace upon
> creation. Such root users have full privileges and it seems that
> these privileges should be controlled through some means (process
> capabilities ?)
>
> Changes [try #2]
>
> - removed struct user_namespace* argument from find_user()
> - added a root_user per user namespace
>
> Signed-off-by: Cedric Le Goater <clg@fr.ibm.com>
> Cc: Andrew Morton <aikpm@osdl.org>
> Cc: Kirill Korotaev <dev@openvz.org>
> Cc: Eric W. Biederman <ebiederm@xmission.com>
> Cc: Herbert Poetzl <herbert@13thfloor.at>
> Cc: Serge E. Hallyn <serue@us.ibm.com>
> Cc: Dave Hansen <haveblue@us.ibm.com>
>
> ---
[... skip ...]
> +struct user_namespace init_user_ns = {
> + .kref = {
> + .refcount = ATOMIC_INIT(2),
<<< please add some comment about why it is initially = 2

> + },
> + .root_user = &root_user,
> +};
> +
> +EXPORT_SYMBOL_GPL(init_user_ns);
>
> /*
> * The uidhash_lock is mostly taken from process context, but it is
> @@ -84,6 +93,111 @@ static inline struct user_struct *uid_ha
> return NULL;
> }

```

>
> +
> +#ifdef CONFIG_USER_NS
> +
> +/*
> + * Clone a new ns copying an original user ns, setting refcount to 1
> + * @old_ns: namespace to clone
> + * Return NULL on error (failure to kmalloc), new ns otherwise
> + */
> +static struct user_namespace *clone_user_ns(struct user_namespace *old_ns)
> +{
> +    struct user_namespace *ns;
> +
> +    ns = kmalloc(sizeof(struct user_namespace), GFP_KERNEL);
> +    if (ns) {
<<<< can you remake this function please so that normal case would go without indentation, i.e.:
if (!ns) {
    error path;
}
normal path

> +    int n;
> +    struct user_struct *new_user;
> +
> +    kref_init(&ns->kref);
> +
> +    for(n = 0; n < UIDHASH_SZ; ++n)
> +        INIT_LIST_HEAD(ns->uidhash_table + n);
> +
> +    /* Insert new root user. */
> +    ns->root_user = alloc_uid(ns, 0);
> +    if (!ns->root_user) {
> +        kfree(ns);
> +        return NULL;
> +    }
> +
> +    /* Reset current->user with a new one */
> +    new_user = alloc_uid(ns, current->uid);
> +    if (!new_user) {
> +        kfree(ns);
> +        return NULL;
> +    }
> +
> +    switch_uid(new_user);
<<<< I see switch_uid in this function, but you don't change nsproxy->user_ns here...
<<<< it is done much later, in sys_unshare()... This looks a bit inconsistent for me...
<<<< But I'm unsure whether it can be fixed... :/

```

```

> + }

> + return ns;
> +}
> +
> +/*
> + * unshare the current process' user namespace.
> + */
> +int unshare_user_ns(unsigned long unshare_flags,
> +    struct user_namespace **new_user)
> +{
> + if (unshare_flags & CLONE_NEWUSER) {
> +   if (!capable(CAP_SYS_ADMIN))
> +     return -EPERM;
<<<< such checks for CAP_SYS_ADMIN mean that we can't use copy_xxx/clone_xxx functions
directly
<<<< from OpenVZ code, since VE creation is done with dropped capabilities already.
<<<< (user level tools decide which capabilities should be granted to VE, so CAP_SYS_ADMIN
<<<< is not normally granted :))
<<<< Can we move capability checks into some more logical place which deals with user, e.g.
sys_unshare()?

> +
> + *new_user = clone_user_ns(current->nsproxy->user_ns);
> + if (!*new_user)
> +   return -ENOMEM;
> +}
> +
> + return 0;
> +}
> +
> +/*
> + * Copy task tsk's user namespace, or clone it if flags specifies
> + * CLONE_NEWUSER. In latter case, changes to the user namespace of
> + * this process won't be seen by parent, and vice versa.
> +*/
> +int copy_user_ns(int flags, struct task_struct *tsk)
> +{
> + struct user_namespace *old_ns = tsk->nsproxy->user_ns;
> + struct user_namespace *new_ns;
> + int err = 0;
> +
> + if (!old_ns)
> +   return 0;
> +
> + get_user_ns(old_ns);
> +
> + if (!(flags & CLONE_NEWUSER))

```

```
> + return 0;  
> +  
> + if (!capable(CAP_SYS_ADMIN)) {  
> + err = -EPERM;  
> + goto out;  
> +}  
<<<< same as above  
  
> +  
> + new_ns = clone_user_ns(old_ns);  
> + if (!new_ns) {  
> + err = -ENOMEM;  
> + goto out;  
> +}  
> + tsk->nsproxy->user_ns = new_ns;  
> +  
> +out:  
> + put_user_ns(old_ns);  
> + return err;  
> +}  
> +
```

....

Kirill
