

On 1/11/07, Serge E. Hallyn <serue@us.ibm.com> wrote:

>  
> That's what's holding me back here - I'm still not sure whether  
> to proceed with a separate implementation, proceed with the  
> current implementation of Paul's containers, or wait for an  
> update from Paul responding to your feedback.

I think that having just a single process-grouping infrastructure in the kernel, rather than separate ones for CPUsets, virtual servers, BeanCounters, ResGroups, etc, is definitely something we should aim for, provided that the requirements for virtual servers aren't too different to those for other users.

I've been thinking about how the `container_clone()` function would need to work. Essentially, the sequence would be something like the following:

1) `do_fork()` or `sys_unshare()` see that some `ns_proxy` bits have changed, and call `container_clone(&nsproxy_subsys)`

In `container_clone()`:

2) Check that the `ns_proxy` container subsystem was bound to some hierarchy; if not, find an unused hierarchy and bind `ns_proxy` subsystem to it. (Or fail with `-EBUSY` if there are no free container subsystems). By the end of this step, `ns_proxy` is bound to hierarchy H.

3) Create a new subcontainer of the current process' container in hierarchy H. This will involve some VFS manipulation, since normally container creation operations are as part of a `mkdir` operation, but shouldn't be too tricky.

4) Move the current process (or possibly the new child process in the case of `do_fork()`) into that container. (This doesn't affect the container mappings of the current or child process in any hierarchies other than H).

So in general if you wanted to combine `ns_proxy` isolation and resource isolation, you'd mount an instance of `containerfs` that bound both the `ns_proxy` subsystem and any resource controllers that you wanted, prior to creating any virtual servers. Alternatively you could have a separate hierarchy for the resource controllers and manually move the new virtual server's root process into the appropriate resource

control container.

Paul

---