

Dave Hansen wrote:

>>This patch is simple, but does not handle SMP scalability very well  
>>(you'll get a lot of cacheline problems when you start actually using  
>>the container structure; the hashing helps a lot there)  
>>  
>>  
>Could you elaborate a bit on this one? What has cacheline problems?  
>  
>

OK, on reflection this probably doesn't belong this early in the series. It only helps speed up lookup by IDs which is only an SMP speed enhancement if you expect a lot of those (eg, when storing the XIDs on another subsystem, such as a filesystem) and unusual things like creating/destroying a lot of vservers quickly. I'll pull it out.

I tried to make it as basic as possible, but my basic problem with the patch you linked to is that it shouldn't be so small that you don't get a decent bunch of the internal API functions listed in description of part 1 of the set  
( [http://vserver.utsi.gen.nz/patches/utsi/2.6.16-rc4-vsi/01-VS\\_erver-Umbrella.diff](http://vserver.utsi.gen.nz/patches/utsi/2.6.16-rc4-vsi/01-VS_erver-Umbrella.diff))

>>and does not provide functions such as looking up a container by ID etc.  
>>  
>>  
>  
>We need something so simple that we probably don't even deal with ids.  
>I believe that Eric claims that we don't really need container \_ids\_.  
>For instance, the filesystem namespaces have no ids, and work just fine.  
>  
>

For actual namespace structures themselves that are virtualising real things, I agree entirely.

But I think the whole point of this patchset (which sadly vger ate for the wider audience, and I still don't know why) is to group tasks and to give userland, and hence the administrator, something tangible with which to group processes with and tack namespace structures onto. I can split out the ID related stuff into another patch, but it would need to go back in before a useful syscall interface can be added. I'll consider it anyway, though.

Note that a given vx\_info structure still might share namespace structures with other vx\_info objects - this is what I was alluding to with the "we haven't made any restrictions on the nature of virtualisation yet" comment. All we've made is the (XID, PID) tuple (although the default for XID=0 is that all processes are visible in one big PID space). The intention is that flags will control how you want your PIDs to work for that vserver.

The only thing that they can't do is share processes - it is a one to many relationship. However, if there can be a parent/child containership relationship on the XIDs themselves, you can still achieve the behaviour of these unusual situations. I'm working on the assumption that if we ever have to migrate trees of XIDs then we can virtualise how they look inside a vserver using flags.

Eric, perhaps you can comment or refer to the earlier post where you made this argument. I tried to follow it but perhaps I missed the jist of one of the messages or one of the most important messages entirely.

>By starting painfully simply, we can build on complexity in bits, and  
>justify it as we go.  
>  
>

This is my aim too! I'll keep chopping and changing.

>>And also because the acronym "vx" makes the API look nice, at least to  
>>mine and Herbert's eyes, then when you go to the network virtualisation  
>>you get "nx\_info", etc. However I'm thinking any of these terms might  
>>also be right:

>>  
>> - "vserver" spelt in full  
>> - family  
>> - container  
>> - jail  
>> - task\_ns (sort for namespace)

>>  
>>Perhaps we can get a ruling from core team on this one, as it's  
>>aesthetics :-).

>>  
>>

>  
>I was in a meeting with a few coworkers, and we were arguing a bit about  
>naming. One person there was a manager-type who didn't have any direct  
>involvement in the project. We asked him which naming was more clear.

>  
>We need to think a bit like that. What is more clear to somebody who  
>has never read the code? (Hint "vx\_" means nothing. :)

>  
>

OK, so let's look at all of the various names that stem from the use of the term "XID" and try to come up with a good naming system for each. I invite the OpenVZ team, Eric and anyone else to put forward their names as well.

Sorry if this is a bit long. Again I invite anyone at all to come forward with a preference or list another set of alternatives. Let's nail this one down, it comes up every time any patchset like this is put forward.

Linux-VServer:

CONFIG\_VSERVER - config option  
typedef unsigned int xid\_t;  
struct vx\_info;  
task\_struct->vx\_info  
task\_struct->xid  
vx\_task\_xid(struct task\*) - get an XID from a task\_struct  
vx\_current\_xid - get XID for current  
vx\_info\_state(struct vx\_info\*, VXS\_FOO) - does vx\_info have state  
create\_vx\_info - creates a new context and "hashes" it  
lookup\_vx\_info - lookup a vx\_info by xid  
get\_vx\_info - increase refcount of a vx\_info  
[...]  
release\_vx\_info - decrease the process count for a vx\_info  
task\_get\_vx\_info - like get\_vx\_info, but by process  
vx\_migrate\_task - join task to a vx\_info  
vxlprintk - debugging printk (for CONFIG\_VSERVER\_DEBUG)  
vxh\_alloc\_vx\_info - history tracing (for CONFIG\_VSERVER\_HISTORY)  
constants:  
  VXS\_FOO - state bits  
  VXF\_FOO - vserver flags (to select features)  
  VXC\_FOO - vserver-specific capabilities  
  VCMD\_get\_version - vserver subcommand names  
  VCI\_VERSION - perhaps the legacy of this one should die.

Using the term "vserver" and ID term "vsid":

CONFIG\_VSERVER - config option  
typedef unsigned int vsid\_t;  
struct vserver  
task\_struct->vserver  
task\_struct->vsid  
vserver\_task\_vsid(struct task\*) - get an VSID from a task\_struct  
vserver\_current\_vsid - get VSID for current

vserver\_state(struct vserver\*, VS\_STATE\_FOO) - does vserver hav...  
 create\_vserver - creates a new context and "hashes" it  
 lookup\_vserver - lookup a vserver by vsid  
 get\_vserver - increase refcount of a vserver  
 [...]

release\_vserver - decrease the process count for a vserver  
 task\_get\_vserver - like get\_vserver, but by process  
 vserver\_migrate\_task - join task to a vserver  
 vserver\_debug - debugging printk (for CONFIG\_VSERVER\_DEBUG)  
 vserver\_hist\_alloc\_vserver - history tracing (for CONFIG\_VSERVER...  
 constants:  
 VS\_STATE\_FOO - state bits  
 VS\_FLAG\_FOO - vserver flags (to select features)  
 VS\_CAP\_FOO - vserver-specific capabilities  
 VS\_CMD\_get\_version - vserver subcommand names  
 VS\_VCI\_VERSION

Using the term "container" and ID term "cid":

CONFIG\_CONTAINERS - config option  
 typedef unsigned int cid\_t;  
 struct container  
 task\_struct->container  
 task\_struct->cid  
 container\_task\_cid(struct task\*) - get an CID from a task\_struct  
 container\_current\_cid - get CID for current  
 container\_state(struct container\*, CONTAINER\_STATE\_FOO) - does c...  
 create\_container - creates a new context and "hashes" it  
 lookup\_container - lookup a container by cid  
 get\_container - increase refcount of a container  
 [...]

release\_container - decrease the process count for a container  
 task\_get\_container - like get\_container, but by process  
 contain\_task - join task to a container  
 container\_debug - debugging printk (for CONFIG\_CONTAINER\_DEBUG)  
 container\_hist\_alloc\_container - history tracing (for CONFIG\_CON...  
 constants:  
 CONTAINER\_STATE\_FOO - state bits  
 CONTAINER\_FLAG\_FOO - container flags (to select features)  
 CONTAINER\_CAP\_FOO - container-specific capabilities  
 CONTAINER\_CMD\_get\_version - container subcommand names  
 CONTAINER\_VCI\_VERSION

Using the term "box" and ID term "boxid":

CONFIG\_BOXES - config option  
 typedef unsigned int boxid\_t;  
 struct box

task\_struct->box  
 task\_struct->boxid  
 box\_task\_boxid(struct task\*) - get an BOXID from a task\_struct  
 box\_current\_boxid - get BOXID for current  
 box\_state(struct box\*, BOX\_STATE\_FOO) - does box have state  
 create\_box - creates a new context and "hashes" it  
 lookup\_box - lookup a box by boxid  
 get\_box - increase refcount of a box  
 [...]

release\_box - decrease the process count for a box  
 task\_get\_box - like get\_box, but by process  
 box\_migrate\_task - join task to a box  
 box\_printk - debugging printk (for CONFIG\_BOXES\_DEBUG)  
 box\_hist\_alloc\_box - history tracing (for CONFIG\_BOXES\_HISTORY)

constants:

- BOX\_STATE\_FOO - state bits
- BOX\_FLAG\_FOO - box flags (to select features)
- BOX\_CAP\_FOO - box-specific capabilities
- BOX\_CMD\_get\_version - box subcommand names
- BOX\_VCI\_VERSION

Using the term "family" and ID term "fid":

CONFIG\_FAMILY - config option  
 typedef unsigned int fid\_t;  
 struct family  
 task\_struct->family  
 task\_struct->fid  
 task\_fid(struct task\*) - get an FID from a task\_struct  
 family\_current\_fid - get FID for current  
 family\_state(struct family\*, FAMILY\_STATE\_FOO) - does family hav...  
 create\_family - creates a new context and "hashes" it  
 lookup\_family - lookup a family by fid  
 get\_family - increase refcount of a family  
 [...]

release\_family - decrease the process count for a family  
 task\_get\_family - like get\_family, but by process  
 family\_adopt\_task - join task to a family  
 family\_printk - debugging printk (for CONFIG\_FAMILY\_DEBUG)  
 family\_hist\_alloc\_family - history tracing (for CONFIG\_FAMILY\_HI...)

constants:

- FAMILY\_STATE\_FOO - state bits
- FAMILY\_FLAG\_FOO - family flags (to select features)
- FAMILY\_CAP\_FOO - family-specific capabilities
- FAMILY\_CMD\_get\_version - family subcommand names
- FAMILY\_VCI\_VERSION

Using the term "task\_ns" and ID term "nsid":

CONFIG\_TASK\_NS - config option  
typedef unsigned int nsid\_t;  
struct task\_ns  
task\_struct->task\_ns  
task\_struct->nsid  
task\_nsid(struct task\*) - get an NSID from a task\_struct  
current\_nsid - get NSID for current  
task\_ns\_state(struct task\_ns\*, TASK\_NS\_STATE\_FOO) - does task\_ns hav...  
create\_task\_ns - creates a new context and "hashes" it  
lookup\_task\_ns - lookup a task\_ns by nsid  
get\_task\_ns - increase refcount of a task\_ns  
[...]  
release\_task\_ns - decrease the process count for a task\_ns  
task\_get\_task\_ns - like get\_task\_ns, but by process  
task\_ns\_migrate\_task - join task to a task\_ns  
task\_ns\_printk - debugging printk (for CONFIG\_TASK\_NS\_DEBUG)  
task\_ns\_hist\_alloc\_task\_ns - history tracing (for CONFIG\_TASK\_NS\_HI...  
constants:  
TASK\_NS\_STATE\_FOO - state bits  
TASK\_NS\_FLAG\_FOO - task\_ns flags (to select features)  
TASK\_NS\_CAP\_FOO - task\_ns-specific capabilities  
TASK\_NS\_CMD\_get\_version - task\_ns subcommand names  
TASK\_NS\_VCI\_VERSION

For the record, I like the term "process family" most. It implies the possibility strict grouping, like last name, as well as allowing heirarchies, but of course in our modern, post-nuclear family age, does not imply a fixed nature of anything ;-).

Happy picking.

Sam.