

---

Subject: Re: [PATCH 1/9] Containers (V9): Basic container framework  
Posted by [Paul Menage](#) on Tue, 01 May 2007 17:46:10 GMT  
[View Forum Message](#) <> [Reply to Message](#)

---

On 5/1/07, Balbir Singh <balbir@linux.vnet.ibm.com> wrote:

> menage@google.com wrote:

> > This patch adds the main containers framework - the container  
> > filesystem, and the basic structures for tracking membership and  
> > associating subsystem state objects to tasks.

>

> [snip]

>

> > +\*\*\* notify\_on\_release is disabled in the current patch set. It may be

> > +\*\*\* reactivated in a future patch in a less-intrusive manner

> > +

>

> Won't this break user space tools for cpusets?

Yes, so it's a must-fix before this gets anywhere near a real distribution.

>

> [snip]

>

> > +See kernel/container.c for more details.

> > +

> > +Subsystems can take/release the container\_mutex via the functions

> > +container\_lock()/container\_unlock(), and can

> > +take/release the callback\_mutex via the functions

> > +container\_lock()/container\_unlock().

> > +

>

> Hmm.. looks like a documentation error. Both mutex's are obtained through

> container\_lock/container\_unlock ?

The second half of that sentence is obsolete and should have been deleted.

>

> > +Accessing a task's container pointer may be done in the following ways:

> > +- while holding container\_mutex

> > +- while holding the task's alloc\_lock (via task\_lock())

> > +- inside an rcu\_read\_lock() section via rcu\_dereference()

> > +

>

> container\_mutex() and task\_lock() can be used for changing the pointer?

No, these are all for read operations. (Actually, this is a bit of documentation that's bit-rotted - there's no longer a per-task "container" pointer). I'll update this.

For write operations, only the container system should be modifying those pointers (under the protection of both container\_mutex and alloc\_lock).

>  
> We needed the equivalent of container\_remove\_file() to be called  
> if container\_add\_file() failed.  
>

Yes, this is some incomplete behaviour that I inherited from cpusets. Needs tidying up.

>  
> Can't we derive the top\_container from containerfs\_root?

Yes, we could for the cost of an extra dereference. Not sure it's a big deal either way.

> > + ssize\_t (\*read) (struct container \*cont, struct cftype \*cft,  
> > + struct file \*file,  
> > + char \_\_user \*buf, size\_t nbytes, loff\_t \*ppos);  
> > + u64 (\*read\_uint) (struct container \*cont, struct cftype \*cft);  
>  
> Is this a new callback, a specialization of the read() callback?

Yes. It's to simplify the common case of reporting a number in a control file. (Not yet well documented :-)

Paul

---