

---

Subject: Re: [PATCH v6 10/10] Disable task moving when using kernel memory accounting

Posted by [Glauber Costa](#) on Mon, 28 Nov 2011 11:00:24 GMT

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

---

On 11/28/2011 02:32 AM, KAMEZAWA Hiroyuki wrote:

> On Fri, 25 Nov 2011 15:38:16 -0200

> Glauber Costa<glommer@parallels.com> wrote:

>  
>> Since this code is still experimental, we are leaving the exact  
>> details of how to move tasks between cgroups when kernel memory  
>> accounting is used as future work.  
>>  
>> For now, we simply disallow movement if there are any pending  
>> accounted memory.  
>>  
>> Signed-off-by: Glauber Costa<glommer@parallels.com>  
>> CC: Hiroyouki Kamezawa<kamezawa.hiroyu@jp.fujitsu.com>  
>> ---  
>> mm/memcontrol.c | 23 ++++++  
>> 1 files changed, 22 insertions(+), 1 deletions(-)  
>>  
>> diff --git a/mm/memcontrol.c b/mm/memcontrol.c  
>> index 2df5d3c..ab7e57b 100644  
>> --- a/mm/memcontrol.c  
>> +++ b/mm/memcontrol.c  
>> @@ -5451,10 +5451,19 @@ static int mem\_cgroup\_can\_attach(struct cgroup\_subsys \*ss,  
>> {  
>> int ret = 0;  
>> struct mem\_cgroup \*mem = mem\_cgroup\_from\_cont(cgroup);  
>> + struct mem\_cgroup \*from = mem\_cgroup\_from\_task(p);  
>> +  
>> + #if defined(CONFIG\_CGROUP\_MEM\_RES\_CTLR\_KMEM) && defined(CONFIG\_INET)  
>> + if (from != mem && !mem\_cgroup\_is\_root(from) &&  
>> + res\_counter\_read\_u64(&from->tcp\_mem.tcp\_memory\_allocated, RES\_USAGE)) {  
>> + printk(KERN\_WARNING "Can't move tasks between cgroups: "  
>> + "Kernel memory held. task: %s\n", p->comm);  
>> + return 1;  
>> + }  
>> + #endif  
>>  
> Hmm, the kernel memory is not guaranteed as being held by the 'task' ?  
>  
> How about  
> "Now, moving task between cgroup is disallowed while the source cgroup  
> contains kmem reference." ?  
>  
> Hmm.. we need to fix this task-move/rmdir issue before production use.

>  
>  
> Thanks,  
> -Kame  
>  
Hi Kame,

Let me tell you the direction I am going wrt task movement: The only reasons I haven't included so far, is that I believe it needs more testing, and as you know, I am right now more interested in getting past the initial barriers for inclusion. I am committed to fix anything that needs to be fixed - stylish or non-stylish before we remove the experimental flag.

So what I intend to do, is to basically

- \* lock the task,
- \* scan through its file descriptors list,
- \* identify which of them are sockets,
- \* cast them to struct sock \*,
- \* see if it has a cgrp associated
- \* see if cgrp == from

At this point we can decrement sockets allocated by 1 in from, and memory\_allocated by sk\_forward\_alloc (increasing by equal quantities in the destination cgroup)

I belive it will work.

---