

---

Subject: Re: [RFC/PATCH] cgroup swap subsystem  
Posted by [KAMEZAWA Hiroyuki](#) on Wed, 05 Mar 2008 06:53:29 GMT  
[View Forum Message](#) <> [Reply to Message](#)

---

On Wed, 05 Mar 2008 14:59:05 +0900  
Daisuke Nishimura <nishimura@mxp.nes.nec.co.jp> wrote:

```
> #ifdef CONFIG_CGROUP_MEM_CONT
> +/*
> + * A page_cgroup page is associated with every page descriptor. The
> + * page_cgroup helps us identify information about the cgroup
> + */
> +struct page_cgroup {
> + struct list_head lru; /* per cgroup LRU list */
> + struct page *page;
> + struct mem_cgroup *mem_cgroup;
> +#ifdef CONFIG_CGROUP_SWAP_LIMIT
> + struct mm_struct *pc_mm;
> +#endif
> + atomic_t ref_cnt; /* Helpful when pages move b/w */
> + /* mapped and cached states */
> + int flags;
> +};
>
```

As first impression, I don't like to increase size of this...but have no alternative idea.

```
> static inline int page_cgroup_locked(struct page *page)
> @@ -664,6 +665,10 @@ retry:
> pc->flags = PAGE_CGROUP_FLAG_ACTIVE;
> if (ctype == MEM_CGROUP_CHARGE_TYPE_CACHE)
> pc->flags |= PAGE_CGROUP_FLAG_CACHE;
> +#ifdef CONFIG_CGROUP_SWAP_LIMIT
> + atomic_inc(&mm->mm_count);
> + pc->pc_mm = mm;
> +#endif
>
```

Strongly Nack to this atomic\_inc().

What happens when tmpfs pages goes to swap ?

```
> if (!page || page_cgroup_assign_new_page_cgroup(page, pc)) {
> /*
> @@ -673,6 +678,9 @@ retry:
>
> +int swap_cgroup_charge(struct page *page,
```

```

> + struct swap_info_struct *si,
> + unsigned long offset)
> +{
> + int ret;
> + struct page_cgroup *pc;
> + struct mm_struct *mm;
> + struct swap_cgroup *swap;
> +
> + BUG_ON(!page);
> +
> + /*
> + * Pages to be swapped out should have been charged by memory cgroup,
> + * but very rarely, pc would be NULL (pc is not reliable without lock,
> + * so I should fix here).
> + * In such cases, we charge the init_mm now.
> + */
> + pc = page_get_page_cgroup(page);
> + if (WARN_ON(!pc))
> + mm = &init_mm;
> + else
> + mm = pc->pc_mm;
> + BUG_ON(!mm);
> +
> + rcu_read_lock();
> + swap = rcu_dereference(mm->swap_cgroup);
> + rcu_read_unlock();
> + BUG_ON(!swap);
Is there no race ?

```

At first look, remembering mm struct is not very good.  
Remembering swap controller itself is better.  
If you go this direction, how about this way ?

```

==
enum {
#ifdef CONFIG_CGROUP_MEM_CONT
    MEMORY_RESOURCE_CONTROLLER,
#endif
#ifdef CONFIG_CGROUP_SWAP_CONT
    SWAP_CONTROLLER,
#endif
    NR_PAGE_CONTROLLER,
}

struct page_cgroup {
    .....
    void* controls[NR_PAGE_CONTROLLER];
    ....
}

```

};  
==

Thanks,  
-Kame

---

Containers mailing list  
Containers@lists.linux-foundation.org  
<https://lists.linux-foundation.org/mailman/listinfo/containers>

---