

---

Subject: Re: [PATCH v2 06/11] memcg: kmem controller infrastructure  
Posted by [KAMEZAWA Hiroyuki](#) on Fri, 10 Aug 2012 17:27:17 GMT  
[View Forum Message](#) <> [Reply to Message](#)

---

(2012/08/09 22:01), Glauber Costa wrote:

> This patch introduces infrastructure for tracking kernel memory pages to  
> a given memcg. This will happen whenever the caller includes the flag  
> \_\_GFP\_KMEMCG flag, and the task belong to a memcg other than the root.  
>  
> In memcontrol.h those functions are wrapped in inline accessors. The  
> idea is to later on, patch those with static branches, so we don't incur  
> any overhead when no mem cgroups with limited kmem are being used.  
>  
> [ v2: improved comments and standardized function names ]  
>  
> Signed-off-by: Glauber Costa <glommer@parallels.com>  
> CC: Christoph Lameter <ccl@linux.com>  
> CC: Pekka Enberg <penberg@cs.helsinki.fi>  
> CC: Michal Hocko <mhocko@suse.cz>  
> CC: Kamezawa Hiroyuki <kamezawa.hiroyu@jp.fujitsu.com>  
> CC: Johannes Weiner <jannes@cmpxchg.org>  
> ---  
> include/linux/memcontrol.h | 79 ++++++  
> mm/memcontrol.c | 185 ++++++  
> 2 files changed, 264 insertions(+)  
>  
> diff --git a/include/linux/memcontrol.h b/include/linux/memcontrol.h  
> index 8d9489f..75b247e 100644  
> --- a/include/linux/memcontrol.h  
> +++ b/include/linux/memcontrol.h  
> @@ -21,6 +21,7 @@  
> #define \_LINUX\_MEMCONTROL\_H  
> #include <linux/cgroup.h>  
> #include <linux/vm\_event\_item.h>  
> +#include <linux/hardirq.h>  
>  
> struct mem\_cgroup;  
> struct page\_cgroup;  
> @@ -399,6 +400,11 @@ struct sock;  
> #ifdef CONFIG\_MEMCG\_KMEM  
> void sock\_update\_memcg(struct sock \*sk);  
> void sock\_release\_memcg(struct sock \*sk);  
> +  
> +#define memcg\_kmem\_on 1  
> +bool \_\_memcg\_kmem\_new\_page(gfp\_t gfp, void \*handle, int order);  
> +void \_\_memcg\_kmem\_commit\_page(struct page \*page, void \*handle, int order);  
> +void \_\_memcg\_kmem\_free\_page(struct page \*page, int order);  
> #else

```

> static inline void sock_update_memcg(struct sock *sk)
> {
> @@ -406,6 +412,79 @@ static inline void sock_update_memcg(struct sock *sk)
> static inline void sock_release_memcg(struct sock *sk)
> {
> }
> +
> +#define memcg_kmem_on 0
> +static inline bool
> +__memcg_kmem_new_page(gfp_t gfp, void *handle, int order)
> +{
> + return false;
> +}
> +
> +static inline void __memcg_kmem_free_page(struct page *page, int order)
> +{
> +}
> +
> +static inline void
> +__memcg_kmem_commit_page(struct page *page, struct mem_cgroup *handle, int order)
> +{
> +}
> +
> #endif /* CONFIG_MEMCG_KMEM */
> +
> +/**
> + * memcg_kmem_new_page: verify if a new kmem allocation is allowed.
> + * @gfp: the gfp allocation flags.
> + * @handle: a pointer to the memcg this was charged against.
> + * @order: allocation order.
> + *
> + * returns true if the memcg where the current task belongs can hold this
> + * allocation.
> + *
> + * We return true automatically if this allocation is not to be accounted to
> + * any memcg.
> + */
> +static __always_inline bool
> +memcg_kmem_new_page(gfp_t gfp, void *handle, int order)
> +{
> + if (!memcg_kmem_on)
> + return true;
> + if (!(gfp & __GFP_KMEMCG) || (gfp & __GFP_NOFAIL))
> + return true;
> + if (in_interrupt() || (!current->mm) || (current->flags & PF_KTHREAD))
> + return true;
> + return __memcg_kmem_new_page(gfp, handle, order);
> +}
> +

```

```

> +/*
> + * memcg_kmem_free_page: uncharge pages from memcg
> + * @page: pointer to struct page being freed
> + * @order: allocation order.
> +
> + * there is no need to specify memcg here, since it is embedded in page_cgroup
> +
> +static __always_inline void
> +memcg_kmem_free_page(struct page *page, int order)
> +{
> + if (memcg_kmem_on)
> + __memcg_kmem_free_page(page, order);
> +
> +
> +/*
> + * memcg_kmem_commit_page: embeds correct memcg in a page
> + * @handle: a pointer to the memcg this was charged against.
> + * @page: pointer to struct page recently allocated
> + * @handle: the memcg structure we charged against
> + * @order: allocation order.
> +
> + * Needs to be called after memcg_kmem_new_page, regardless of success or
> + * failure of the allocation. if @page is NULL, this function will revert the
> + * charges. Otherwise, it will commit the memcg given by @handle to the
> + * corresponding page_cgroup.
> +
> +static __always_inline void
> +memcg_kmem_commit_page(struct page *page, struct mem_cgroup *handle, int order)
> +{
> + if (memcg_kmem_on)
> + __memcg_kmem_commit_page(page, handle, order);
> +

```

Doesn't this 2 functions has no short-cuts ?

if (memcg\_kmem\_on && handle) ?

Maybe free() needs to access page\_cgroup...

```

> #endif /* _LINUX_MEMCONTROL_H */
>
> diff --git a/mm/memcontrol.c b/mm/memcontrol.c
> index 54e93de..e9824c1 100644
> --- a/mm/memcontrol.c
> +++ b/mm/memcontrol.c
> @@ -10,6 +10,10 @@

```

```

> * Copyright (C) 2009 Nokia Corporation
> * Author: Kirill A. Shutemov
> *
> + * Kernel Memory Controller
> + * Copyright (C) 2012 Parallels Inc. and Google Inc.
> + * Authors: Glauber Costa and Suleiman Souhlal
> +
>   * This program is free software; you can redistribute it and/or modify
>   * it under the terms of the GNU General Public License as published by
>   * the Free Software Foundation; either version 2 of the License, or
> @@ -434,6 +438,9 @@ struct mem_cgroup *mem_cgroup_from_css(struct
cgroup_subsys_state *s)
> #include <net/ip.h>
>
> static bool mem_cgroup_is_root(struct mem_cgroup *memcg);
> +static int memcg_charge_kmem(struct mem_cgroup *memcg, gfp_t gfp, s64 delta);
> +static void memcg_uncharge_kmem(struct mem_cgroup *memcg, s64 delta);
> +
> void sock_update_memcg(struct sock *sk)
> {
>   if (mem_cgroup_sockets_enabled) {
> @@ -488,6 +495,118 @@ struct cg_proto *tcp_proto_cgroup(struct mem_cgroup *memcg)
>   }
>   EXPORT_SYMBOL(tcp_proto_cgroup);
> #endif /* CONFIG_INET */
> +
> +static inline bool memcg_kmem_enabled(struct mem_cgroup *memcg)
> +{
> + return !mem_cgroup_disabled() && !mem_cgroup_is_root(memcg) &&
> + memcg->kmem_accounted;
> +}
> +
> +/*
> + * We need to verify if the allocation against current->mm->owner's memcg is
> + * possible for the given order. But the page is not allocated yet, so we'll
> + * need a further commit step to do the final arrangements.
> +
> + * It is possible for the task to switch cgroups in this mean time, so at
> + * commit time, we can't rely on task conversion any longer. We'll then use
> + * the handle argument to return to the caller which cgroup we should commit
> + * against
> +
> + * Returning true means the allocation is possible.
> +*/
> +bool __memcg_kmem_new_page(gfp_t gfp, void *_handle, int order)
> +{
> + struct mem_cgroup *memcg;
> + struct mem_cgroup **handle = (struct mem_cgroup **)_handle;

```

```

> + bool ret = true;
> + size_t size;
> + struct task_struct *p;
> +
> + *handle = NULL;
> + rCU_read_lock();
> + p = rCU_dereference(current->mm->owner);
> + memcg = mem_cgroup_from_task(p);
> + if (!memcg_kmem_enabled(memcg))
> +     goto out;
> +
> + mem_cgroup_get(memcg);
> +

```

This mem\_cgroup\_get() will be a potential performance problem.  
 Don't you have good idea to avoid accessing atomic counter here ?  
 I think some kind of percpu counter or a feature to disable "move task"  
 will be a help.

```

> + size = PAGE_SIZE << order;
> + ret = memcg_charge_kmem(memcg, gfp, size) == 0;
> + if (!ret) {
> +     mem_cgroup_put(memcg);
> +     goto out;
> +
> +
> + *handle = memcg;
> +out:
> + rCU_read_unlock();
> + return ret;
> +
> +EXPORT_SYMBOL(__memcg_kmem_new_page);
> +
> +void __memcg_kmem_commit_page(struct page *page, void *handle, int order)
> +{
> +    struct page_cgroup *pc;
> +    struct mem_cgroup *memcg = handle;
> +
> +    if (!memcg)
> +        return;
> +
> +    WARN_ON(mem_cgroup_is_root(memcg));
> +    /* The page allocation must have failed. Revert */
> +    if (!page) {
> +        size_t size = PAGE_SIZE << order;
> +
> +        memcg_uncharge_kmem(memcg, size);

```

```

> + mem_cgroup_put(memcg);
> +
> + return;
> +
> +
> + pc = lookup_page_cgroup(page);
> + lock_page_cgroup(pc);
> + pc->mem_cgroup = memcg;
> + SetPageCgroupUsed(pc);
> + unlock_page_cgroup(pc);
> +
> +
> +void __memcg_kmem_free_page(struct page *page, int order)
> +{
> + struct mem_cgroup *memcg;
> + size_t size;
> + struct page_cgroup *pc;
> +
> + if (mem_cgroup_disabled())
> + return;
> +
> + pc = lookup_page_cgroup(page);
> + lock_page_cgroup(pc);
> + memcg = pc->mem_cgroup;
> + pc->mem_cgroup = NULL;

```

shouldn't this happen after checking "Used" bit ?  
Ah, BTW, why do you need to clear pc->memcg ?

```

> + if (!PageCgroupUsed(pc)) {
> + unlock_page_cgroup(pc);
> + return;
> +
> + ClearPageCgroupUsed(pc);
> + unlock_page_cgroup(pc);
> +
> + /*
> + * Checking if kmem accounted is enabled won't work for uncharge, since
> + * it is possible that the user enabled kmem tracking, allocated, and
> + * then disabled it again.
> +
> + * We trust if there is a memcg associated with the page, it is a valid
> + * allocation
> +
> + */
> + if (!memcg)
> + return;
> +
> + WARN_ON(mem_cgroup_is_root(memcg));
> + size = (1 << order) << PAGE_SHIFT;

```

```
> + memcg_uncharge_kmem(memcg, size);
> + mem_cgroup_put(memcg);
```

Why do we need ref-counting here ? kmem res\_counter cannot work as reference ?

```
> +}
> +EXPORT_SYMBOL(__memcg_kmem_free_page);
> #endif /* CONFIG_MEMCG_KMEM */
>
> #if defined(CONFIG_INET) && defined(CONFIG_MEMCG_KMEM)
> @@ -5759,3 +5878,69 @@ static int __init enable_swap_account(char *s)
>   __setup("swapaccount=", enable_swap_account);
>
> #endif
> +
> +#ifdef CONFIG_MEMCG_KMEM
> +int memcg_charge_kmem(struct mem_cgroup *memcg, gfp_t gfp, s64 delta)
> +{
```

What does 'delta' means ?

```
> + struct res_counter *fail_res;
> + struct mem_cgroup *_memcg;
> + int ret;
> + bool may_oom;
> + bool nofail = false;
> +
> + may_oom = (gfp & __GFP_WAIT) && (gfp & __GFP_FS) &&
> +   !(gfp & __GFP_NORETRY);
> +
> + ret = 0;
> +
> + if (!memcg)
> +   return ret;
> +
> + _memcg = memcg;
> + ret = __mem_cgroup_try_charge(NULL, gfp, delta / PAGE_SIZE,
> +   &_memcg, may_oom);
> +
> + if (ret == -EINTR) {
> +   nofail = true;
> + /*
> +   * __mem_cgroup_try_charge() choiced to bypass to root due to
> +   * OOM kill or fatal signal. Since our only options are to
> +   * either fail the allocation or charge it to this cgroup, do
```

```
> + * it as a temporary condition. But we can't fail. From a
> + * kmem/slab perspective, the cache has already been selected,
> + * by mem_cgroup_get_kmem_cache(), so it is too late to change
> + * our minds
> + */
> + res_counter_charge_nofail(&memcg->res, delta, &fail_res);
> + if (do_swap_account)
> + res_counter_charge_nofail(&memcg->memsw, delta,
> +     &fail_res);
> + ret = 0;
```

Hm, you returns 0 and this charge may never be uncharged....right ?

Thanks,  
-Kame

---