

---

Subject: [PATCH v4 11/25] consider a memcg parameter in kmem\_create\_cache

Posted by [Glauber Costa](#) on Mon, 18 Jun 2012 10:28:04 GMT

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

---

Allow a memcg parameter to be passed during cache creation.

When the slab allocator is being used, it will only merge  
caches that belong to the same memcg.

Default function is created as a wrapper, passing NULL  
to the memcg version. We only merge caches that belong  
to the same memcg.

>From the memcontrol.c side, 3 helper functions are created:

- 1) memcg\_css\_id: because slab needs a unique cache name  
for sysfs. Since this is visible, but not the canonical  
location for slab data, the cache name is not used, the  
css\_id should suffice.
- 2) mem\_cgroup\_register\_cache: is responsible for assigning  
a unique index to each cache, and other general purpose  
setup. The index is only assigned for the root caches. All  
others are assigned index == -1.
- 3) mem\_cgroup\_release\_cache: can be called from the root cache  
destruction, and will release the index for  
other caches.

We can't assign indexes until the basic slab is up and running  
this is because the ida subsystem will itself call slab functions  
such as kmalloc a couple of times. Because of that, we have  
a late\_initcall that scan all caches and register them after the  
kernel is booted up. Only caches registered after that receive  
their index right away.

This index mechanism was developed by Suleiman Souhlal.  
Changed to a idr/ida based approach based on suggestion  
from Kamezawa.

Signed-off-by: Glauber Costa <[glommer@parallels.com](mailto:glommer@parallels.com)>

CC: Christoph Lameter <[cl@linux.com](mailto:cl@linux.com)>

CC: Pekka Enberg <[penberg@cs.helsinki.fi](mailto:penberg@cs.helsinki.fi)>

CC: Michal Hocko <[mhocko@suse.cz](mailto:mhocko@suse.cz)>

CC: Kamezawa Hiroyuki <[kamezawa.hiroyu@jp.fujitsu.com](mailto:kamezawa.hiroyu@jp.fujitsu.com)>

CC: Johannes Weiner <[hannes@cmpxchg.org](mailto:hannes@cmpxchg.org)>

CC: Suleiman Souhlal <[suleiman@google.com](mailto:suleiman@google.com)>

---

include/linux/memcontrol.h | 14 ++++++++-----

```
include/linux/slab.h      | 10 ++++++++
mm/memcontrol.c          | 25 ++++++=====
mm/slab.h                | 22 ++++++=====
mm/slab_common.c         | 36 ++++++=====
mm/slub.c                | 16 ++++++=====
6 files changed, 109 insertions(+), 14 deletions(-)
```

```
diff --git a/include/linux/memcontrol.h b/include/linux/memcontrol.h
index f94efd2..99e14b9 100644
--- a/include/linux/memcontrol.h
+++ b/include/linux/memcontrol.h
@@ -26,6 +26,7 @@ struct mem_cgroup;
struct page_cgroup;
struct page;
struct mm_struct;
+struct kmem_cache;

/* Stats that can be updated by kernel.*/
enum mem_cgroup_page_stat_item {
@@ -440,7 +441,20 @@ struct sock;
#endif CONFIG_CGROUP_MEM_RES_CTLR_KMEM
void sock_update_memcg(struct sock *sk);
void sock_release_memcg(struct sock *sk);
+int memcg_css_id(struct mem_cgroup *memcg);
+void mem_cgroup_register_cache(struct mem_cgroup *memcg,
+       struct kmem_cache *s);
+void mem_cgroup_release_cache(struct kmem_cache *cachep);
#else
+static inline void mem_cgroup_register_cache(struct mem_cgroup *memcg,
+       struct kmem_cache *s)
+{
+}
+
+static inline void mem_cgroup_release_cache(struct kmem_cache *cachep)
+{
+}
+
 static inline void sock_update_memcg(struct sock *sk)
{
}

diff --git a/include/linux/slab.h b/include/linux/slab.h
index e007cc5..d347616 100644
--- a/include/linux/slab.h
+++ b/include/linux/slab.h
@@ -116,6 +116,7 @@ struct kmem_cache {
};

#endif
```

```

+struct mem_cgroup;
/*
 * struct kmem_cache related prototypes
 */
@@ -125,6 +126,9 @@ int slab_is_available(void);
struct kmem_cache *kmem_cache_create(const char *, size_t, size_t,
 unsigned long,
 void (*)(void *));
+struct kmem_cache *
+kmem_cache_create_memcg(struct mem_cgroup *, const char *, size_t, size_t,
+ unsigned long, void (*)(void *));
void kmem_cache_destroy(struct kmem_cache *);
int kmem_cache_shrink(struct kmem_cache *);
void kmem_cache_free(struct kmem_cache *, void *);
@@ -338,6 +342,12 @@ extern void *__kmalloc_track_caller(size_t, gfp_t, unsigned long);
 __kmalloc(size, flags)
#endif /* DEBUG_SLAB */

+ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
+#define MAX_KMEM_CACHE_TYPES 400
+else
+#define MAX_KMEM_CACHE_TYPES 0
+endif /* CONFIG_CGROUP_MEM_RES_CTLR_KMEM */
+
#ifdef CONFIG_NUMA
/*
 * kmalloc_node_track_caller is a special version of kmalloc_node that
diff --git a/mm/memcontrol.c b/mm/memcontrol.c
index 04ac774..cb57c5b 100644
--- a/mm/memcontrol.c
+++ b/mm/memcontrol.c
@@ -323,6 +323,11 @@ struct mem_cgroup {
#endif
};

+int memcg_css_id(struct mem_cgroup *memcg)
+{
+ return css_id(&memcg->css);
+}
+
/* Stuffs for move charges at task migration. */
/*
 * Types of charges to be moved. "move_charge_at_immitgrate" is treated as a
@@ -461,6 +466,25 @@ struct cg_proto *tcp_proto_cgroup(struct mem_cgroup *memcg)
}
EXPORT_SYMBOL(tcp_proto_cgroup);
#endif /* CONFIG_INET */
+

```

```

+struct ida cache_types;
+
+void mem_cgroup_register_cache(struct mem_cgroup *memcg,
+      struct kmem_cache *cachep)
+{
+ int id = -1;
+
+ if (!memcg)
+ id = ida_simple_get(&cache_types, 0, MAX_KMEM_CACHE_TYPES,
+ GFP_KERNEL);
+ cachep->memcg_params.id = id;
+}
+
+void mem_cgroup_release_cache(struct kmem_cache *cachep)
+{
+ if (cachep->memcg_params.id != -1)
+ ida_simple_remove(&cache_types, cachep->memcg_params.id);
+}
#endif /* CONFIG_CGROUP_MEM_RES_CTLR_KMEM */

static void drain_all_stock_async(struct mem_cgroup *memcg);
@@ -5053,6 +5077,7 @@ mem_cgroup_create(struct cgroup *cont)
#endif CONFIG_CGROUP_MEM_RES_CTLR_KMEM
    WARN_ON(cgroup_add_cftypes(&mem_cgroup_subsys,
        kmem_cgroup_files));
+ ida_init(&cache_types);
#endif

    if (mem_cgroup_soft_limit_tree_init())
diff --git a/mm/slab.h b/mm/slab.h
index 19e17c7..1781580 100644
--- a/mm/slab.h
+++ b/mm/slab.h
@@ -39,16 +39,28 @@ unsigned long calculate_alignment(unsigned long flags,
int __kmem_cache_create(struct kmem_cache *s);
int __kmem_cache_initcall(void);

+struct mem_cgroup;
#ifndef CONFIG_SLUB
-struct kmem_cache *__kmem_cache_alias(const char *name, size_t size,
- size_t align, unsigned long flags, void (*ctor)(void *));
+struct kmem_cache *
+__kmem_cache_alias(struct mem_cgroup *memcg, const char *name, size_t size,
+ size_t align, unsigned long flags, void (*ctor)(void *));
#else
-static inline struct kmem_cache *__kmem_cache_alias(const char *name, size_t size,
- size_t align, unsigned long flags, void (*ctor)(void__))
+static inline struct kmem_cache *

```

```

+__kmem_cache_alias(struct mem_cgroup *memcg, const char *name, size_t size,
+    size_t align, unsigned long flags, void (*ctor)(void *))
{ return NULL; }
#endif

-
int __kmem_cache_shutdown(struct kmem_cache *);

+static inline bool cache_match_memcg(struct kmem_cache *cachep,
+    struct mem_cgroup *memcg)
+{
+#ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
+    return cachep->memcg_params.memcg == memcg;
+#endif
+    return true;
+}
+
+void __init memcg_slab_register_all(void);
#endif

diff --git a/mm/slab_common.c b/mm/slab_common.c
index 15db694ac..42f226d 100644
--- a/mm/slab_common.c
+++ b/mm/slab_common.c
@@ -16,6 +16,7 @@ 
#include <asm/cacheflush.h>
#include <asm/tlbflush.h>
#include <asm/page.h>
+#include <linux/memcontrol.h>

#include "slab.h"

@@ -77,8 +78,9 @@ unsigned long calculate_alignment(unsigned long flags,
 * as davem.
 */
-struct kmem_cache *kmem_cache_create(const char *name, size_t size, size_t align,
-    unsigned long flags, void (*ctor)(void *))
+struct kmem_cache *
+kmem_cache_create_memcg(struct mem_cgroup *memcg, const char *name, size_t size,
+    size_t align, unsigned long flags, void (*ctor)(void *))
{
    struct kmem_cache *s = NULL;
    char *n;
@@ -118,7 +120,7 @@ struct kmem_cache *kmem_cache_create(const char *name, size_t size,
size_t align
    continue;
}

```

```

- if (!strcmp(s->name, name)) {
+ if (cache_match_memcg(s, memcg) && !strcmp(s->name, name)) {
    printk(KERN_ERR "kmem_cache_create(%s): Cache name"
        " already exists.\n",
        name);
@@ -131,7 +133,7 @@ struct kmem_cache *kmem_cache_create(const char *name, size_t size,
size_t align
    WARN_ON(strchr(name, ' ')); /* It confuses parsers */
#endif

- s = __kmem_cache_alias(name, size, align, flags, ctor);
+ s = __kmem_cache_alias(memcg, name, size, align, flags, ctor);
if (s)
    goto oops;

@@ -152,11 +154,17 @@ struct kmem_cache *kmem_cache_create(const char *name, size_t
size, size_t align
    s->flags = flags;
    s->align = calculate_alignment(flags, align, size);

+#ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
+ s->memcg_params.memcg = memcg;
+#endif
+
r = __kmem_cache_create(s);

if (!r) {
    s->refcount = 1;
    list_add(&s->list, &slab_caches);
+ if (slab_state >= FULL)
+     mem_cgroup_register_cache(memcg, s);
}
else {
    kmem_cache_free(kmem_cache, s);
@@ -173,6 +181,12 @@ oops:

    return s;
}
+
+struct kmem_cache *kmem_cache_create(const char *name, size_t size, size_t align,
+ unsigned long flags, void (*ctor)(void *))
+{
+ return kmem_cache_create_memcg(NULL, name, size, align, flags, ctor);
+}
EXPORT_SYMBOL(kmem_cache_create);

void kmem_cache_destroy(struct kmem_cache *s)
@@ -185,6 +199,7 @@ void kmem_cache_destroy(struct kmem_cache *s)

```

```

if (s->flags & SLAB_DESTROY_BY_RCU)
    rCU_barrier();

+ mem_cgroup_release_cache(s);
    kfree(s->name);
    kmem_cache_free(kmem_cache, s);
} else {
@@ -205,6 +220,17 @@ int slab_is_available(void)

static int __init kmem_cache_initcall(void)
{
- return __kmem_cache_initcall();
+ int r = __kmem_cache_initcall();
+#ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
+ struct kmem_cache *s;
+
+ if (r)
+     return r;
+ mutex_lock(&slab_mutex);
+ list_for_each_entry(s, &slab_caches, list)
+     mem_cgroup_register_cache(NULL, s);
+ mutex_unlock(&slab_mutex);
#endif
+ return r;
}
__initcall(kmem_cache_initcall);
diff --git a/mm/slub.c b/mm/slub.c
index d9d4d5a..ca4b8e0 100644
--- a/mm/slub.c
+++ b/mm/slub.c
@@ -31,6 +31,7 @@ 
#include <linux/fault-inject.h>
#include <linux/stacktrace.h>
#include <linux/prefetch.h>
+#include <linux/memcontrol.h>

#include <trace/events/kmem.h>

@@ -3831,7 +3832,7 @@ static int slab_unmergeable(struct kmem_cache *s)
    return 0;
}

-static struct kmem_cache *find_mergeable(size_t size,
+static struct kmem_cache *find_mergeable(struct mem_cgroup *memcg, size_t size,
    size_t align, unsigned long flags, const char *name,
    void (*ctor)(void *))
{
@@ -3867,17 +3868,20 @@ static struct kmem_cache *find_mergeable(size_t size,

```

```

if (s->size - size >= sizeof(void *))
    continue;

+ if (!cache_match_memcg(s, memcg))
+ continue;
    return s;
}
return NULL;
}

-struct kmem_cache *__kmem_cache_alias(const char *name, size_t size,
- size_t align, unsigned long flags, void (*ctor)(void *))
+struct kmem_cache *
+__kmem_cache_alias(struct mem_cgroup *memcg, const char *name, size_t size,
+ size_t align, unsigned long flags, void (*ctor)(void *))
{
    struct kmem_cache *s;

- s = find_mergeable(size, align, flags, name, ctor);
+ s = find_mergeable(memcg, size, align, flags, name, ctor);
    if (s) {
        s->refcount++;
        /*
@@@ -5207,6 +5211,10 @@ static char *create_unique_id(struct kmem_cache *s)
        if (p != name + 1)
            *p++ = '-';
        p += sprintf(p, "%07d", s->size);
+#ifdef CONFIG_CGROUP_MEM_RES_CTLR_KMEM
+ if (s->memcg_params.memcg)
+     p += sprintf(p, "-%08d", memcg_css_id(s->memcg_params.memcg));
#endif
        BUG_ON(p > name + ID_STR_LENGTH - 1);
        return name;
    }
--
```

1.7.10.2

---