Subject: [PATCH v4 15/19] memcg/sl[au]b: shrink dead caches Posted by Glauber Costa on Fri, 12 Oct 2012 13:41:09 GMT

View Forum Message <> Reply to Message

In the slub allocator, when the last object of a page goes away, we don't necessarily free it - there is not necessarily a test for empty page in any slab_free path.

This means that when we destroy a memcg cache that happened to be empty, those caches may take a lot of time to go away: removing the memcg reference won't destroy them - because there are pending references, and the empty pages will stay there, until a shrinker is called upon for any reason.

This patch marks all memcg caches as dead. kmem_cache_shrink is called for the ones who are not yet dead - this will force internal cache reorganization, and then all references to empty pages will be removed.

An unlikely branch is used to make sure this case does not affect performance in the usual slab_free path.

The slab allocator has a time based reaper that would eventually get rid of the objects, but we can also call it explicitly, since dead caches are not a likely event.

```
[ v2: also call verify_dead for the slab ]
[ v3: use delayed_work to avoid calling verify_dead at every free]
Signed-off-by: Glauber Costa <glommer@parallels.com>
CC: Christoph Lameter <cl@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 <hannes@cmpxchg.org>
CC: Suleiman Souhlal <suleiman@google.com>
CC: Tejun Heo <tj@kernel.org>
include/linux/slab.h | 2 +-
mm/memcontrol.c
                    2 files changed, 23 insertions(+), 7 deletions(-)
diff --git a/include/linux/slab.h b/include/linux/slab.h
index e17d348..9ed1a98 100644
--- a/include/linux/slab.h
+++ b/include/linux/slab.h
@ @ -214,7 +214,7 @ @ struct memcg_cache_params {
  struct kmem_cache *cachep;
  bool dead:
```

```
atomic_t nr_pages;

    struct work struct destroy;

+ struct delayed_work destroy;
 };
};
};
diff --git a/mm/memcontrol.c b/mm/memcontrol.c
index 8cf8b4d..b52e6b9 100644
--- a/mm/memcontrol.c
+++ b/mm/memcontrol.c
@ @ -2955,11 +2955,27 @ @ static void kmem_cache_destroy_work_func(struct work_struct *w)
 struct kmem_cache *cachep;
 struct memcg_cache_params *p;
+ struct delayed_work *dw = to_delayed_work(w);
- p = container of(w, struct memcg cache params, destroy);
+ p = container_of(dw, struct memcg_cache_params, destroy);
 cachep = p->cachep;
- if (!atomic read(&cachep->memcg params->nr pages))
+ /*
+ * If we get down to 0 after shrink, we could delete right away.
+ * However, memcg_release_pages() already puts us back in the workqueue
+ * in that case. If we proceed deleting, we'll get a dangling
+ * reference, and removing the object from the workqueue in that case is
+ * unnecessary complication. We are not a fast path.
+ * If we aren't down to zero, we'll schedule a slow worker and try again
+ if (atomic read(&cachep->memcg params->nr pages) != 0) {
+ kmem cache shrink(cachep);
+ if (atomic_read(&cachep->memcg_params->nr_pages) == 0)
+ return;
+ /* Once per minute should be good enough. */
+ schedule delayed work(&cachep->memcq params->destroy, 60 * HZ);
+ } else
 kmem cache destroy(cachep);
static DECLARE WORK(kmem cache destroy work, kmem cache destroy work func);
@@ -2973,7 +2989,7 @@ void mem_cgroup_destroy_cache(struct kmem_cache *cachep)
 * We have to defer the actual destroying to a workqueue, because
 * we might currently be in a context that cannot sleep.
 */
- schedule_work(&cachep->memcg_params->destroy);
+ schedule_delayed_work(&cachep->memcg_params->destroy, 0);
}
```