

---

Subject: Re: [PATCH 8/8] Per-container pages reclamation

Posted by [Balbir Singh](#) on Fri, 01 Jun 2007 07:02:33 GMT

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

---

Andrew Morton wrote:

> On Wed, 30 May 2007 19:42:26 +0400

> Pavel Emelianov <xemul@openvz.org> wrote:

>

>> Implement try\_to\_free\_pages\_in\_container() to free the

>> pages in container that has run out of memory.

>>

>> The scan\_control->isolate\_pages() function is set to

>> isolate\_pages\_in\_container() that isolates the container

>> pages only. The exported \_\_isolate\_lru\_page() call

>> makes things look simpler than in the previous version.

>>

>> Includes fix from Balbir Singh <balbir@in.ibm.com>

>>

>> }

>>

>> +void container\_rss\_move\_lists(struct page \*pg, bool active)

>> +{

>> + struct rss\_container \*rss;

>> + struct page\_container \*pc;

>> +

>> + if (!page\_mapped(pg))

>> + return;

>> +

>> + pc = page\_container(pg);

>> + if (pc == NULL)

>> + return;

>> +

>> + rss = pc->cnt;

>> +

>> + spin\_lock(&rss->res.lock);

>> + if (active)

>> + list\_move(&pc->list, &rss->active\_list);

>> + else

>> + list\_move(&pc->list, &rss->inactive\_list);

>> + spin\_unlock(&rss->res.lock);

>> +}

>

> This is an interesting-looking function. Please document it?

>

Will do. This function is called when we want to move a page in the LRU list. This could happen when a page is activated or when reclaim finds that a particular page cannot be reclaimed right now.

> I'm inferring that the rss container has an active and inactive list and  
> that this basically follows the same operation as the traditional per-zone  
> lists?  
>

Yes, correct.

> Would I be correct in guessing that pages which are on the  
> per-rss-container lists are also eligible for reclaim off the traditional  
> page LRUs? If so, how does that work? When a page gets freed off the  
> per-zone LRUs does it also get removed from the per-rss\_container LRU? But  
> how can this be right? Pages can get taken off the LRU and freed at  
> interrupt time, and this code isn't interrupt-safe.  
>

Yes, before a page is reclaimed from the global LRU list, we go through  
page\_remove\_rmap() in try\_to\_unmap(). Pages are removed from the container  
LRU before they are reclaimed.

> I note that this lock is not irq-safe, whereas the lru locks are irq-safe.  
> So we don't perform the rotate\_reclaimable\_page() operation within the RSS  
> container? I think we could do so. I wonder if this was considered.  
>

The lock needs to be interrupt safe.

> A description of how all this code works would help a lot.

```
>
>> +static unsigned long isolate_container_pages(unsigned long nr_to_scan,
>> + struct list_head *src, struct list_head *dst,
>> + unsigned long *scanned, struct zone *zone, int mode)
>> +{
>> + unsigned long nr_taken = 0;
>> + struct page *page;
>> + struct page_container *pc;
>> + unsigned long scan;
>> + LIST_HEAD(pc_list);
>> +
>> + for (scan = 0; scan < nr_to_scan && !list_empty(src); scan++) {
>> + pc = list_entry(src->prev, struct page_container, list);
>> + page = pc->page;
>> + if (page_zone(page) != zone)
>> + continue;
```

>  
> That page\_zone() check is interesting. What's going on here?  
>

> I'm suspecting that we have a problem here: if there are a lot of pages on

> \*src which are in the wrong zone, we can suffer reclaim distress leading to  
> omm-killings, or excessive CPU consumption?  
>

We discussed this on lkml. Basically, now for every zone we try to reclaim pages from the container, it increases CPU utilization if we choose the wrong zone to reclaim from, but it provides the following benefit

Code reuse (shrink\_zone\* is reused along with helper functions). I am not sure if Pavel had any other benefits in mind like benefits on a NUMA box.

```
>> + for_each_online_node(node) {  
>> + #ifdef CONFIG_HIGHMEM  
>> + zones = NODE_DATA(node)->node_zonelist[ZONE_HIGHMEM].zones;  
>> + if (do_try_to_free_pages(zones, sc.gfp_mask, &sc))  
>> + return 1;  
>> + #endif  
>> + zones = NODE_DATA(node)->node_zonelist[ZONE_NORMAL].zones;  
>> + if (do_try_to_free_pages(zones, sc.gfp_mask, &sc))  
>> + return 1;  
>> + }  
>  
> Definitely need to handle ZONE_DMA32 and ZONE_DMA (some architectures put  
> all memory into ZONE_DMA (or they used to))  
>
```

node\_zonelist[ZONE\_NORMAL].zones should contain ZONE\_DMA and ZONE\_DMA32 right?

--

Warm Regards,  
Balbir Singh  
Linux Technology Center  
IBM, ISTL

---