
Subject: [PATCH 2/6] BC: bean counters core (API)
Posted by [dev](#) on Wed, 23 Aug 2006 11:00:37 GMT

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

Core functionality and interfaces of BC:
find/create bean counter, initialization,
charge/uncharge of resource, core objects' declarations.

Basic structures:

bc_resource_parm - resource description
bean counter - set of resources, id, lock

Signed-off-by: Pavel Emelianov <xemul@sw.ru>

Signed-off-by: Kirill Korotaev <dev@sw.ru>

```
include/bc/beancounter.h | 140 ++++++  
init/main.c | 4  
kernel/Makefile | 1  
kernel/bc/Makefile | 7 +  
kernel/bc/beancounter.c | 292 ++++++  
5 files changed, 444 insertions(+)
```

```
--- /dev/null 2006-07-18 14:52:43.075228448 +0400  
+++ ./include/bc/beancounter.h 2006-08-21 13:14:01.000000000 +0400  
@@ -0,0 +1,140 @@  
+/*  
+ * include/bc/beancounter.h  
+ *  
+ * Copyright (C) 2006 OpenVZ. SWsoft Inc  
+ *  
+ */  
+  
+#ifndef _LINUX_BEANCOUNTER_H  
+#define _LINUX_BEANCOUNTER_H  
+  
+/*  
+ * Resource list.  
+ */  
+  
+#define BC_RESOURCES 0  
+  
+struct bc_resource_parm {  
+ unsigned long barrier; /* A barrier over which resource allocations  
+ * are failed gracefully. e.g. if the amount  
+ * of consumed memory is over the barrier
```

```

+   * further sbrk() or mmap() calls fail, the
+   * existing processes are not killed.
+   */
+ unsigned long limit; /* hard resource limit */
+ unsigned long held; /* consumed resources */
+ unsigned long maxheld; /* maximum amount of consumed resources */
+ unsigned long minheld; /* minimum amount of consumed resources */
+ unsigned long failcnt; /* count of failed charges */
+};

+
+/*
+ * Kernel internal part.
+ */
+
+ifdef __KERNEL__
+
+#include <linux/config.h>
+#include <linux/spinlock.h>
+#include <linux/list.h>
+#include <asm/atomic.h>
+
#define BC_MAXVALUE LONG_MAX
+
+/*
+ * Resource management structures
+ * Serialization issues:
+ *   beancounter list management is protected via bc_hash_lock
+ *   task pointers are set only for current task and only once
+ *   refcount is managed atomically
+ *   value and limit comparison and change are protected by per-bc spinlock
+ */
+
+struct beancounter
+{
+ atomic_t bc_refcount;
+ spinlock_t bc_lock;
+ uid_t bc_id;
+ struct hlist_node hash;
+
+ /* resources statistics and settings */
+ struct bc_resource_parm bc_parms[BC_RESOURCES];
+};
+
+enum severity { BC_BARRIER, BC_LIMIT, BC_FORCE };
+
+/* Flags passed to beancounter_findcreate() */
#define BC_ALLOC 0x01 /* May allocate new one */
#define BC_ALLOC_ATOMIC 0x02 /* Allocate with GFP_ATOMIC */

```

```

+
+">#define BC_HASH_SIZE 256
+
+ifdef CONFIG_BEANCOUNTERS
+extern struct hlist_head bc_hash[];
+extern spinlock_t bc_hash_lock;
+
+/*
+ * This function tunes minheld and maxheld values for a given
+ * resource when held value changes
+ */
+static inline void bc_adjust_held_minmax(struct beancounter *bc,
+ int resource)
+{
+ if (bc->bc_parms[resource].maxheld < bc->bc_parms[resource].held)
+ bc->bc_parms[resource].maxheld = bc->bc_parms[resource].held;
+ if (bc->bc_parms[resource].minheld > bc->bc_parms[resource].held)
+ bc->bc_parms[resource].minheld = bc->bc_parms[resource].held;
+}
+
+void bc_print_resource_warning(struct beancounter *bc, int res,
+ char *str, unsigned long val, unsigned long held);
+void bc_print_id(struct beancounter *bc, char *str, int size);
+
+int bc_charge_locked(struct beancounter *bc,
+ int res, unsigned long val, enum severity strict);
+int bc_charge(struct beancounter *bc,
+ int res, unsigned long val, enum severity strict);
+
+void bc_uncharge_locked(struct beancounter *bc,
+ int res, unsigned long val);
+void bc_uncharge(struct beancounter *bc,
+ int res, unsigned long val);
+
+struct beancounter *beancounter_findcreate(uid_t id, int mask);
+
+static inline struct beancounter *get_beancounter(struct beancounter *bc)
+{
+ atomic_inc(&bc->bc_refcount);
+ return bc;
+}
+
+void __put_beancounter(struct beancounter *bc);
+static inline void put_beancounter(struct beancounter *bc)
+{
+ __put_beancounter(bc);
+}
+

```

```

+void bc_init_early(void);
+void bc_init_late(void);
+void bc_init_proc(void);
+
+extern struct beancounter init_bc;
+extern const char *bc_rnames[];
+
+/*#else /* CONFIG_BEANCOUNTERS */
+
+#define beancounter_findcreate(id, f) (NULL)
+#define get_beancounter(bc) (NULL)
+#define put_beancounter(bc) do { } while (0)
+#define bc_charge_locked(bc, r, v, s) (0)
+#define bc_charge(bc, r, v) (0)
+#define bc_uncharge_locked(bc, r, v) do { } while (0)
+#define bc_uncharge(bc, r, v) do { } while (0)
+#define bc_init_early() do { } while (0)
+#define bc_init_late() do { } while (0)
+#define bc_init_proc() do { } while (0)
+
+/*#endif /* CONFIG_BEANCOUNTERS */
+/*#endif /* __KERNEL__ */
+
+/*#endif /* _LINUX_BEANCOUNTER_H */
--- ./init/main.c.ve1 2006-08-21 12:25:13.000000000 +0400
+++ ./init/main.c 2006-08-21 12:45:32.000000000 +0400
@@ -52,6 +52,8 @@
#include <linux/debug_locks.h>
#include <linux/lockdep.h>

+#include <bc/beancounter.h>
+
#include <asm/io.h>
#include <asm/bugs.h>
#include <asm/setup.h>
@@ -494,6 +496,7 @@ asmlinkage void __init start_kernel(void
early_boot_irqs_off();
early_init_irq_lock_class();

+ bc_init_early();
/*
 * Interrupts are still disabled. Do necessary setups, then
 * enable them
@@ -587,6 +590,7 @@ asmlinkage void __init start_kernel(void
#endif
fork_init(num_physpages);
proc_caches_init();
+ bc_init_late();

```

```

buffer_init();
unnamed_dev_init();
key_init();
--- ./kernel/Makefile.ve1 2006-08-21 12:25:14.000000000 +0400
+++ ./kernel/Makefile 2006-08-21 12:25:27.000000000 +0400
@@ @ -12,6 +12,7 @@ obj-y    = sched.o fork.o exec_domain.o

obj-$(CONFIG_STACKTRACE) += stacktrace.o
obj-y += time/
+obj-y += bc/
obj-$(CONFIG_DEBUG_MUTEXES) += mutex-debug.o
obj-$(CONFIG_LOCKDEP) += lockdep.o
ifeq ($(CONFIG_PROC_FS),y)
--- /dev/null 2006-07-18 14:52:43.075228448 +0400
+++ ./kernel/bc/Makefile 2006-08-21 12:25:27.000000000 +0400
@@ @ -0,0 +1,7 @@
+#
+## Beancounters (BC)
+#
+## Copyright (C) 2006 OpenVZ. SWsoft Inc
+#
+
+obj-$(CONFIG_BEANCOUNTERS) += beancounter.o
--- /dev/null 2006-07-18 14:52:43.075228448 +0400
+++ ./kernel/bc/beancounter.c 2006-08-21 13:13:11.000000000 +0400
@@ @ -0,0 +1,292 @@
+/*
+ * kernel/bc/beancounter.c
+ *
+ * Copyright (C) 2006 OpenVZ. SWsoft Inc
+ * Original code by (C) 1998 Alan Cox
+ * 1998-2000 Andrey Savochkin <saw@saw.sw.com.sg>
+ */
+
+##include <linux/slab.h>
+##include <linux/module.h>
+
+##include <bc/beancounter.h>
+
+static kmem_cache_t *bc_cachep;
+static struct beancounter default_beancounter;
+
+static void init_beancounter_struct(struct beancounter *bc, uid_t id);
+
+struct beancounter init_bc;
+
+const char *bc_rnames[] = {
+};

```

```

+
+">#define bc_hash_fun(x) (((x) >> 8) ^ (x)) & (BC_HASH_SIZE - 1))
+
+struct hlist_head bc_hash[BC_HASH_SIZE];
+spinlock_t bc_hash_lock;
+
+EXPORT_SYMBOL(bc_hash);
+EXPORT_SYMBOL(bc_hash_lock);
+
+/*
+ * Per resource bean计数。资源是与其luid绑定的。
+ * 资源结构本身既被进程标记，又被正在充电的资源标记（一个socket不想在
+ * 例如irq时间搜索东西）。引用计数器保持东西在手头。
+ *
+ * 在用户创建资源，杀死所有进程，然后启动新进程的情况下，正确处理
+ * 这样。引用计数器将意味着旧条目仍然存在，并与资源绑定。
+ */
+
+struct beancounter *beancounter_findcreate(uid_t uid, int mask)
+{
+    struct beancounter *new_bc, *bc;
+    unsigned long flags;
+    struct hlist_head *slot;
+    struct hlist_node *pos;
+
+    slot = &bc_hash[bc_hash_fun(uid)];
+    new_bc = NULL;
+
+retry:
+    spin_lock_irqsave(&bc_hash_lock, flags);
+    hlist_for_each_entry(bc, pos, slot, hash)
+        if (bc->bc_id == uid)
+            break;
+
+    if (pos != NULL) {
+        get(beancounter(bc));
+        spin_unlock_irqrestore(&bc_hash_lock, flags);
+
+        if (new_bc != NULL)
+            kmem_cache_free(bc_cachep, new_bc);
+        return bc;
+    }
+
+    if (!(mask & BC_ALLOC))
+        goto out_unlock;

```

```

+
+ if (new_bc != NULL)
+ goto out_install;
+
+ spin_unlock_irqrestore(&bc_hash_lock, flags);
+
+ new_bc = kmem_cache_alloc(bc_cachep,
+ mask & BC_ALLOC_ATOMIC ? GFP_ATOMIC : GFP_KERNEL);
+ if (new_bc == NULL)
+ goto out;
+
+ memcpy(new_bc, &default_beancounter, sizeof(*new_bc));
+ init_beancounter_struct(new_bc, uid);
+ goto retry;
+
+out_install:
+ hlist_add_head(&new_bc->hash, slot);
+out_unlock:
+ spin_unlock_irqrestore(&bc_hash_lock, flags);
+out:
+ return new_bc;
+}
+
+void bc_print_id(struct beancounter *bc, char *str, int size)
+{
+ snprintf(str, size, "%u", bc->bc_id);
+}
+
+void bc_print_resource_warning(struct beancounter *bc, int res,
+ char *str, unsigned long val, unsigned long held)
+{
+ char uid[64];
+
+ bc_print_id(bc, uid, sizeof(uid));
+ printk(KERN_WARNING "BC %s %s warning: %s "
+ "(held %lu, fails %lu, val %lu)\n",
+ uid, bc_rnames[res], str,
+ (res < BC_RESOURCES ? bc->bc_parms[res].held : held),
+ (res < BC_RESOURCES ? bc->bc_parms[res].failcnt : 0),
+ val);
+}
+
+static inline void verify_held(struct beancounter *bc)
+{
+ int i;
+
+ for (i = 0; i < BC_RESOURCES; i++)
+ if (bc->bc_parms[i].held != 0)

```

```

+ bc_print_resource_warning(bc, i,
+   "resource is held on put", 0, 0);
+}
+
+void __put_beancounter(struct beancounter *bc)
+{
+ unsigned long flags;
+
+ /* equivalent to atomic_dec_and_lock_irqsave() */
+ local_irq_save(flags);
+ if (unlikely(!atomic_dec_and_lock(&bc->bc_refcount, &bc_hash_lock))) {
+ local_irq_restore(flags);
+ if (unlikely(atomic_read(&bc->bc_refcount) < 0))
+ printk(KERN_ERR "BC: Bad refcount: bc=%p, "
+ "luid=%d, ref=%d\n",
+ bc, bc->bc_id,
+ atomic_read(&bc->bc_refcount));
+ return;
+ }
+
+ BUG_ON(bc == &init_bc);
+ verify_held(bc);
+ hlist_del(&bc->hash);
+ spin_unlock_irqrestore(&bc_hash_lock, flags);
+ kmem_cache_free(bc_cachep, bc);
+}
+
+EXPORT_SYMBOL(__put_beancounter);
+
+/*
+ * Generic resource charging stuff
+ */
+
+/* called with bc->bc_lock held and interrupts disabled */
+int bc_charge_locked(struct beancounter *bc, int resource, unsigned long val,
+ enum severity strict)
+{
+ /*
+ * bc_value <= BC_MAXVALUE, value <= BC_MAXVALUE, and only one addition
+ * at the moment is possible so an overflow is impossible.
+ */
+ bc->bc_parms[resource].held += val;
+
+ switch (strict) {
+ case BC_BARRIER:
+ if (bc->bc_parms[resource].held >
+ bc->bc_parms[resource].barrier)
+ break;

```

```

+ /* fallthrough */
+ case BC_LIMIT:
+ if (bc->bc_parms[resource].held >
+ bc->bc_parms[resource].limit)
+ break;
+ /* fallthrough */
+ case BC_FORCE:
+ bc_adjust_held_minmax(bc, resource);
+ return 0;
+ default:
+ BUG();
+ }
+
+ bc->bc_parms[resource].failcnt++;
+ bc->bc_parms[resource].held -= val;
+ return -ENOMEM;
+}
+EXPORT_SYMBOL(bc_charge_locked);
+
+int bc_charge(struct beancounter *bc, int resource, unsigned long val,
+ enum severity strict)
+{
+ int retval;
+ unsigned long flags;
+
+ BUG_ON(val > BC_MAXVALUE);
+
+ spin_lock_irqsave(&bc->bc_lock, flags);
+ retval = bc_charge_locked(bc, resource, val, strict);
+ spin_unlock_irqrestore(&bc->bc_lock, flags);
+ return retval;
+}
+EXPORT_SYMBOL(bc_charge);
+
+/* called with bc->bc_lock held and interrupts disabled */
+void bc_uncharge_locked(struct beancounter *bc, int resource, unsigned long val)
+{
+ if (unlikely(bc->bc_parms[resource].held < val)) {
+ bc_print_resource_warning(bc, resource,
+ "uncharging too much", val, 0);
+ val = bc->bc_parms[resource].held;
+ }
+
+ bc->bc_parms[resource].held -= val;
+ bc_adjust_held_minmax(bc, resource);
+}
+EXPORT_SYMBOL(bc_uncharge_locked);
+

```

```

+void bc_uncharge(struct beancounter *bc, int resource, unsigned long val)
+{
+    unsigned long flags;
+
+    spin_lock_irqsave(&bc->bc_lock, flags);
+    bc_uncharge_locked(bc, resource, val);
+    spin_unlock_irqrestore(&bc->bc_lock, flags);
+}
+EXPORT_SYMBOL(bc_uncharge);
+
+/*
+ * Initialization
+ *
+ * struct beancounter contains
+ * - limits and other configuration settings
+ * - structural fields: lists, spinlocks and so on.
+ *
+ * Before these parts are initialized, the structure should be memset
+ * to 0 or copied from a known clean structure. That takes care of a lot
+ * of fields not initialized explicitly.
+ */
+
+static void init_beancounter_struct(struct beancounter *bc, uid_t id)
+{
+    atomic_set(&bc->bc_refcount, 1);
+    spin_lock_init(&bc->bc_lock);
+    bc->bc_id = id;
+}
+
+static void init_beancounter_nolimits(struct beancounter *bc)
+{
+    int k;
+
+    for (k = 0; k < BC_RESOURCES; k++) {
+        bc->bc_parms[k].limit = BC_MAXVALUE;
+        bc->bc_parms[k].barrier = BC_MAXVALUE;
+    }
+}
+
+static void init_beancounter_syslimits(struct beancounter *bc)
+{
+    int k;
+
+    for (k = 0; k < BC_RESOURCES; k++)
+        bc->bc_parms[k].barrier = bc->bc_parms[k].limit;
+}
+
+void __init bc_init_early(void)

```

```

+{
+ struct beancounter *bc;
+ struct hlist_head *slot;
+
+ bc = &init_bc;
+
+ memset(bc, 0, sizeof(*bc));
+ init_beancounter_nolimits(bc);
+ init_beancounter_struct(bc, 0);
+
+ spin_lock_init(&bc_hash_lock);
+ slot = &bc_hash[bc_hash_fun(bc->bc_id)];
+ hlist_add_head(&bc->hash, slot);
+}
+
+void __init bc_init_late(void)
+{
+ struct beancounter *bc;
+
+ bc_cachep = kmem_cache_create("beancounters",
+ sizeof(struct beancounter),
+ 0, SLAB_HWCACHE_ALIGN, NULL, NULL);
+ if (bc_cachep == NULL)
+ panic("Can't create bc caches\n");
+
+ bc = &default_beancounter;
+ memset(bc, 0, sizeof(default_beancounter));
+ init_beancounter_syslimits(bc);
+ init_beancounter_struct(bc, 0);
+}

```
