
Subject: [PATCH 3/17] BC: beancounters core (API)

Posted by [dev](#) on Tue, 05 Sep 2006 15:17:35 GMT

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

Core functionality and interfaces of BC:

find/create beancounter, initialization,
charge/uncharge of resource, core objects' declarations.

Basic structures:

bc_resource_parm - resource description
beancounter - 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 | 155 ++++++  
include/linux/types.h   | 16 ++  
init/main.c            |  4  
kernel/Makefile        |  1  
kernel/bc/Makefile     |  7 +  
kernel/bc/beancounter.c | 263 ++++++  
6 files changed, 446 insertions(+)
```

```
--- ./include/bc/beancounter.h.bccore 2006-09-05 12:06:35.000000000 +0400
```

```
+++ ./include/bc/beancounter.h 2006-09-05 12:15:57.000000000 +0400
```

```
@@ @ -0,0 +1,155 @@
```

```
+/*  
+ * 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.
+ */
+
+/* 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;
+ bcid_t bc_id;
+ struct hlist_node hash;
+
+ /* resources statistics and settings */
+ struct bc_resource_parm bc_parms[BC_RESOURCES];
+};
+
+enum bc_severity { BC_BARRIER, BC_LIMIT, BC_FORCE };
+
+/* Flags passed to beancounter_findcreate() */
+#define BC_LOOKUP 0x00
+#define BC_ALLOC 0x01 /* may allocate new one */
+#define BC_ALLOC_ATOMIC 0x02 /* when BC_ALLOC is set causes
+   * GFP_ATOMIC allocation

```

```

+    */
+
+ifdef CONFIG_BEANCOUNTERS
+
+/*
+ * These functions tune minheld and maxheld values for a given
+ * resource when held value changes
+ */
+static inline void bc_adjust_maxheld(struct beancounter *bc, int resource)
+{
+ struct bc_resource_parm *parm;
+
+ parm = &bc->bc_parms[resource];
+ if (parm->maxheld < parm->held)
+     parm->maxheld = parm->held;
+}
+
+static inline void bc_adjust_minheld(struct beancounter *bc, int resource)
+{
+ struct bc_resource_parm *parm;
+
+ parm = &bc->bc_parms[resource];
+ if (parm->minheld > parm->held)
+     parm->minheld = parm->held;
+}
+
+int __must_check bc_charge_locked(struct beancounter *bc,
+ int res, unsigned long val, enum bc_severity strict);
+int __must_check bc_charge(struct beancounter *bc,
+ int res, unsigned long val, enum bc_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(bcid_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);
+
+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[];
+
+/* CONFIG_BEANCOUNTERS */
+
+#define beancounter_findcreate(id, f) (NULL)
+#define get_beancounter(bc) (NULL)
+#define put_beancounter(bc) do { } while (0)
+
+static inline __must_check int bc_charge_locked(struct beancounter *bc,
+ int res, unsigned long val, enum bc_severity strict)
+{
+ return 0;
+}
+
+static inline __must_check int bc_charge(struct beancounter *bc,
+ int res, unsigned long val, enum bc_severity strict)
+{
+ return 0;
+}
+
+static inline void bc_uncharge_locked(struct beancounter *bc, int res,
+ unsigned long val)
+{
+}
+
+static inline void bc_uncharge(struct beancounter *bc, int res,
+ unsigned long val)
+{
+}
+
+#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 */
--- ./include/linux/types.h.bccore 2006-09-05 11:47:33.000000000 +0400
+++ ./include/linux/types.h 2006-09-05 12:06:35.000000000 +0400
@@ -40,6 +40,21 @@ typedef __kernel_gid32_t gid_t;
typedef __kernel_uid16_t uid16_t;
typedef __kernel_gid16_t gid16_t;

+/*
+ * Type of beancounter id (CONFIG_BEANCOUNTERS)
+ *

```

```

+ * The ancient Unix implementations of this kind of resource management and
+ * security are built around setluid() which sets a uid value that cannot
+ * be changed again and is normally used for security purposes. That
+ * happened to be a uid_t and in simple setups at login uid = luid = euid
+ * would be the norm.
+ *
+ * Thus the Linux one happens to be a uid_t. It could be something else but
+ * for the "container per user" model whatever a container is must be able
+ * to hold all possible uid_t values. Alan Cox.
+ */
+typedef uid_t bcid_t;
+
#ifndef CONFIG_UID16
/* This is defined by include/asm-{arch}/posix_types.h */
typedef __kernel_old_uid_t old_uid_t;
@@ -52,6 +67,7 @@ typedef __kernel_old_gid_t old_gid_t;
#else
typedef __kernel_uid_t uid_t;
typedef __kernel_gid_t gid_t;
+typedef __kernel_uid_t bcid_t;
#endif /* __KERNEL__ */

#if defined(__GNUC__) && !defined(__STRICT_ANSI__)
--- ./init/main.c.bccore 2006-09-05 11:47:33.000000000 +0400
+++ ./init/main.c 2006-09-05 12:06:35.000000000 +0400
@@ -50,6 +50,8 @@
#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>
@@ -493,6 +495,7 @@ 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
@@ -585,6 +588,7 @@ 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.bccore 2006-09-05 11:47:33.000000000 +0400
+++ ./kernel/Makefile 2006-09-05 12:09:53.000000000 +0400
@@ @ -12,6 +12,7 @@ obj-y    = sched.o fork.o exec_domain.o

obj-$(CONFIG_STACKTRACE) += stacktrace.o
obj-y += time/
+obj-$(CONFIG_BEANCOUNTERS) += bc/
obj-$(CONFIG_DEBUG_MUTEXES) += mutex-debug.o
obj-$(CONFIG_LOCKDEP) += lockdep.o
ifeq ($(CONFIG_PROC_FS),y)
--- ./kernel/bc/Makefile.bccore 2006-09-05 12:06:35.000000000 +0400
+++ ./kernel/bc/Makefile 2006-09-05 12:10:05.000000000 +0400
@@ @ -0,0 +1,7 @@
@#
+#
+## Bean counters (BC)
+#
+## Copyright (C) 2006 OpenVZ. SWsoft Inc
+#
+
+obj-y += beancounter.o
--- ./kernel/bc/beancounter.c.bccore 2006-09-05 12:06:35.000000000 +0400
+++ ./kernel/bc/beancounter.c 2006-09-05 12:16:50.000000000 +0400
@@ @ -0,0 +1,263 @@
+/*
+ * 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 <linux/hash.h>
+
+##include <bc/beancounter.h>
+
+static kmem_cache_t *bc_cachep;
+static struct beancounter default_beancounter;
+
+static void init_beancounter_struct(struct beancounter *bc, bcid_t id);
+
+struct beancounter init_bc;
+
+const char *bc_rnames[] = {
+};

```

```

+
+">#define BC_HASH_BITS 8
+[#define BC_HASH_SIZE (1 << BC_HASH_BITS)
+
+static struct hlist_head bc_hash[BC_HASH_SIZE];
+static spinlock_t bc_hash_lock;
+[#define bc_hash_fn(bcid) (hash_long(bcid, BC_HASH_BITS))
+
+/*
+ * Per resource bean计数。资源与他们的bc id绑定。
+ * 资源结构本身既被进程和正在充电的资源（一个socket不想在 irq 时搜索）标记。
+ * 参考计数器保持东西在手头。
+ *
+ * 在一个用户创建资源，杀死所有他的进程，然后启动新的时候，这样处理是正确的。参考计数器
+ * 将意味着旧的条目仍然存在，并且资源与之绑定。
+ */
+
+struct beancounter *beancounter_findcreate(bcid_t id, int mask)
+{
+    struct beancounter *new_bc, *bc;
+    unsigned long flags;
+    struct hlist_head *slot;
+    struct hlist_node *pos;
+
+    slot = &bc_hash[bc_hash_fn(id)];
+    new_bc = NULL;
+
+retry:
+    spin_lock_irqsave(&bc_hash_lock, flags);
+    hlist_for_each_entry(bc, pos, slot, hash)
+        if (bc->bc_id == id)
+            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 (new_bc != NULL)
+        goto out_install;
+

```

```

+ spin_unlock_irqrestore(&bc_hash_lock, flags);
+
+ if (!(mask & BC_ALLOC))
+ goto out;
+
+ new_bc = kmem_cache_alloc(bc_cachep,
+ mask & BC_ALLOC_ATOMIC ? GFP_ATOMIC : GFP_KERNEL);
+ if (new_bc == NULL)
+ goto out;
+
+ *new_bc = default_beancounter;
+ init_beancounter_struct(new_bc, id);
+ goto retry;
+
+out_install:
+ hlist_add_head(&new_bc->hash, slot);
+ spin_unlock_irqrestore(&bc_hash_lock, flags);
+out:
+ return new_bc;
+}
+
+void put_beancounter(struct beancounter *bc)
+{
+ int i;
+ unsigned long flags;
+
+ if (!atomic_dec_and_lock_irqsave(&bc->bc_refcount,
+ &bc_hash_lock, flags))
+ return;
+
+ BUG_ON(bc == &init_bc);
+
+ for (i = 0; i < BC_RESOURCES; i++)
+ if (bc->bc_parms[i].held != 0)
+ printk("BC: %d has %lu of %s held on put", bc->bc_id,
+ bc->bc_parms[i].held, bc_rnames[i]);
+
+ hlist_del(&bc->hash);
+ spin_unlock_irqrestore(&bc_hash_lock, flags);
+
+ kmem_cache_free(bc_cachep, bc);
+}
+
+EXPORT_SYMBOL_GPL(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 bc_severity strict)
+{
+ unsigned long new_held;
+
+ /*
+ * bc_value <= BC_MAXVALUE, value <= BC_MAXVALUE, and only one addition
+ * at the moment is possible so an overflow is impossible.
+ */
+ new_held = 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->bc_parms[resource].held = new_held;
+ bc_adjust_maxheld(bc, resource);
+ return 0;
+
+ default:
+ BUG();
+ }
+
+ bc->bc_parms[resource].failcnt++;
+ return -ENOMEM;
+}
+EXPORT_SYMBOL_GPL(bc_charge_locked);
+
+int bc_charge(struct beancounter *bc, int resource, unsigned long val,
+ enum bc_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_GPL(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)) {
+   printk("BC: overuncharging bc %d %s: val %lu, holds %lu\n",
+     bc->bc_id, bc_rnames[resource], val,
+     bc->bc_parms[resource].held);
+   val = bc->bc_parms[resource].held;
+ }
+
+ bc->bc_parms[resource].held -= val;
+ bc_adjust_minheld(bc, resource);
+}
+EXPORT_SYMBOL_GPL(bc_uncharge_locked);
+
+void bc_uncharge(struct beancounter *bc, int resource, unsigned long val)
+{
+ unsigned long flags;
+
+ BUG_ON(val > BC_MAXVALUE);
+
+ spin_lock_irqsave(&bc->bc_lock, flags);
+ bc_uncharge_locked(bc, resource, val);
+ spin_unlock_irqrestore(&bc->bc_lock, flags);
+}
+EXPORT_SYMBOL_GPL(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, bcid_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;
+
+ init_beancounter_nolimits(bc);
+ init_beancounter_struct(bc, 0);
+
+ spin_lock_init(&bc_hash_lock);
+ slot = &bc_hash[bc_hash_fn(bc->bc_id)];
+ hlist_add_head(&bc->hash, slot);
+}
+
+void __init bc_init_late(void)
+{
+ struct beancounter *bc;
+
+ bc_cachep = kmalloc_cache_create("beancounters",
+ sizeof(struct beancounter), 0,
+ SLAB_HWCACHE_ALIGN | SLAB_PANIC, NULL, NULL);
+
+ bc = &default_beancounter;
+ init_beancounter_syslimits(bc);
+ init_beancounter_struct(bc, 0);
+}
```
