
Subject: [PATCH 8/10] BC: private pages (mappings) accounting (core)

Posted by [dev](#) on Thu, 05 Oct 2006 15:54:25 GMT

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

Accounting of private pages.

Accounting is performed on mmap per vma. When mm changes its BC, vmas are not recharged.

On mlock and mprotect vma status may change from BC point of view, so recharging is done accordingly.

See definition of private pages at

http://wiki.openvz.org/User_pages_accounting

Signed-off-by: Pavel Emelianov <xemul@openvz.org>

Signed-off-by: Kirill Korotaev <dev@openvz.org>

```
include/bc/vmpages.h |  90 ++++++=====
include/linux/mm.h   |   3 +
include/linux/sched.h|   3 +
kernel/bc/vmpages.c | 138 ++++++=====
4 files changed, 234 insertions(+)
```

```
--- /dev/null 2006-07-18 14:52:43.075228448 +0400
+++ ./include/bc/vmpages.h 2006-10-05 12:19:59.000000000 +0400
@@ -0,0 +1,90 @@
+/*
+ * include/bc/vmpages.h
+ *
+ * Copyright (C) 2006 OpenVZ SWsoft Inc
+ *
+ */
+
+ifndef __BC_VMPAGES_H_
+define __BC_VMPAGES_H_
+
+#include <bc/beancounter.h>
+
+struct vm_area_struct;
+struct mm_struct;
+struct file;
+
#define BC_NOCHARGE 0
#define BC_UNCHARGE 1
#define BC_CHARGE 2
+
```

```

+ifdef CONFIG_BEANCOUNTERS
#define __vma_set_bc(vma, bc) do { (vma)->vma_bc = bc_get(bc); } while (0)
#define vma_set_bc(vma) __vma_set_bc(vma, (vma)->vm_mm->mm_bc)
#define vma_copy_bc(vma) __vma_set_bc(vma, (vma)->vma_bc)
#define vma_release_bc(vma) do { bc_put((vma)->vma_bc); } while (0)
+
#define mm_init_beancounter(mm) do { \
+ struct beancounter *bc; \
+ bc = get_exec_bc(); \
+ (mm)->mm_bc = bc_get(bc); \
+ } while (0)
#define mm_free_beancounter(mm) do { bc_put(mm->mm_bc); } while (0)
+
+int __must_check bc_need_memory_recharge(struct vm_area_struct *vma,
+ struct file *new_file, unsigned long new_flags);
+
+int __must_check __bc_memory_charge(struct mm_struct *mm, unsigned long len,
+ int severity);
+int __must_check bc_memory_charge(struct mm_struct *mm, unsigned long len,
+ struct file *file, unsigned long flags, int severity);
+int __must_check bc_vma_charge(struct vm_area_struct *vma);
+
+void __bc_memory_uncharge(struct mm_struct *mm, unsigned long len);
+void bc_memory_uncharge(struct mm_struct *mm, unsigned long len,
+ struct file *file, unsigned long flags);
+void bc_vma_uncharge(struct vm_area_struct *vma);
+
#define bc_equal(bc1, bc2) (bc1 == bc2)
#else
static inline
+int __must_check bc_need_memory_recharge(struct vm_area_struct *vma,
+ struct file *new_file, unsigned long new_flags)
+{
+ return BC_NOCHARGE;
+}
static inline int __must_check __bc_memory_charge(struct mm_struct *mm,
+ unsigned long len, int severity)
+{
+ return 0;
+}
static inline int __must_check bc_memory_charge(struct mm_struct *mm,
+ unsigned long len, struct file *file, unsigned long flags,
+ int severity)
+{
+ return 0;
+}
static inline int __must_check bc_vma_charge(struct vm_area_struct *vma)
+{

```

```

+ return 0;
+}
+static inline void __bc_memory_uncharge(struct mm_struct *mm, unsigned long len)
+{
+}
+static inline void bc_memory_uncharge(struct mm_struct *mm, unsigned long len,
+ struct file *file, unsigned long flags)
+{
+}
+static inline void bc_vma_uncharge(struct vm_area_struct *vma)
+{
+}
+
+#define mm_init_beancounter(mm) do { } while (0)
+#define mm_free_beancounter(mm) do { } while (0)
+#define __vma_set_bc(vma, bc) do { } while (0)
+#define vma_set_bc(vma) do { } while (0)
+#define vma_copy_bc(vma) do { } while (0)
+#define vma_release_bc(vma) do { } while (0)
+#define bc_equal(bc1, bc2) 1
#endif
#endif
--- ./include/linux/mm.h.bc_vmpages_core 2006-10-05 12:12:24.000000000 +0400
+++ ./include/linux/mm.h 2006-10-05 12:16:03.000000000 +0400
@@ -112,6 +112,9 @@ struct vm_area_struct {
#ifndef CONFIG_NUMA
    struct mempolicy *vm_policy; /* NUMA policy for the VMA */
#endif
#ifndef CONFIG_BEANCOUNTERS
    struct beancounter *vma_bc;
#endif
};

/*
--- ./include/linux/sched.h.bc_vmpages_core 2006-10-05 12:11:44.000000000 +0400
+++ ./include/linux/sched.h 2006-10-05 12:16:03.000000000 +0400
@@ -369,6 +369,9 @@ struct mm_struct {
/* aio bits */
    rwlock_t ioctx_list_lock;
    struct ioctx *ioctx_list;
#ifndef CONFIG_BEANCOUNTERS
    struct beancounter *mm_bc;
#endif
};

struct sighand_struct {
--- /dev/null 2006-07-18 14:52:43.075228448 +0400
+++ ./kernel/bc/vmpages.c 2006-10-05 12:20:08.000000000 +0400

```

```

@@ -0,0 +1,138 @@
+/*
+ * kernel/bc/vmpages.c
+ *
+ * Copyright (C) 2006 OpenVZ SWsoft Inc
+ *
+ */
+
+[#include <linux/mm.h>
+
+[#include <bc/beancounter.h>
+[#include <bc/vmpages.h>
+
+[#define BC_PRIVVMPAGES_BARRIER BC_MAXVALUE
+[#define BC_PRIVVMPAGES_LIMIT BC_MAXVALUE
+
+/*
+ * Core routines
+ */
+
+/*
+ * bc_vma_private checks whether VMA (file, flags) is private
+ * from BC point of view. private VMAs are charged when they are mapped
+ * thus preventing system from resource exhausting when pages from these VMAs
+ * are touched.
+ */
+static inline int bc_vma_private(struct file *file, unsigned long flags)
+{
+    return (flags & VM_LOCKED) ||
+        ((flags & VM_WRITE) && (file == NULL || !(flags & VM_SHARED)));
+}
+
+/*
+ * Accounting is performed in pages (not in Kbytes)
+ */
+static inline int do_memory_charge(struct beancounter *bc,
+    unsigned long len, int severity)
+{
+    return bc_charge(bc, BC_PRIVVMPAGES, len >> PAGE_SHIFT, severity);
+}
+
+static inline void do_memory_uncharge(struct beancounter *bc, unsigned long len)
+{
+    bc_uncharge(bc, BC_PRIVVMPAGES, len >> PAGE_SHIFT);
+}
+
+/*
+ * API calls

```

```

+ */
+
+int __bc_memory_charge(struct mm_struct *mm, unsigned long len, int severity)
+{
+    return do_memory_charge(mm->mm_bc, len, severity);
+}
+
+int bc_memory_charge(struct mm_struct *mm, unsigned long len,
+    struct file *file, unsigned long flags, int severity)
+{
+    int ret;
+
+    ret = 0;
+    if (bc_vma_private(file, flags))
+        ret = do_memory_charge(mm->mm_bc, len, severity);
+    return ret;
+}
+
+int bc_vma_charge(struct vm_area_struct *vma)
+{
+    int ret;
+
+    ret = (bc_vma_private(vma->vm_file, vma->vm_flags) ?
+        do_memory_charge(vma->vm_mm->mm_bc,
+            vma->vm_end - vma->vm_start, BC_BARRIER) : 0);
+    if (ret == 0)
+        vma_set_bc(vma);
+    return ret;
+}
+
+void __bc_memory_uncharge(struct mm_struct *mm, unsigned long len)
+{
+    do_memory_uncharge(mm->mm_bc, len);
+}
+
+void bc_memory_uncharge(struct mm_struct *mm, unsigned long len,
+    struct file *file, unsigned long flags)
+{
+    if (bc_vma_private(file, flags))
+        do_memory_uncharge(mm->mm_bc, len);
+}
+
+void bc_vma_uncharge(struct vm_area_struct *vma)
+{
+    if (bc_vma_private(vma->vm_file, vma->vm_flags))
+        do_memory_uncharge(vma->vma_bc, vma->vm_end - vma->vm_start);
+    vma_release_bc(vma);
+}

```

```

+
+
+int bc_need_memory_recharge(struct vm_area_struct *vma, struct file *new_file,
+ unsigned long new_flags)
+{
+ if (bc_vma_private(vma->vm_file, vma->vm_flags)) {
+ if (bc_vma_private(new_file, new_flags))
+ return BC_NOCHARGE;
+
+ /* private -> non-private */
+ return BC_UNCHARGE;
+ } else {
+ if (!bc_vma_private(new_file, new_flags))
+ return BC_NOCHARGE;
+
+ /* non-private -> private */
+ return BC_CHARGE;
+ }
+}
+
+/*
+ * Generic resource info
+ */
+
+static int bc_privvm_init(struct beancounter *bc, gfp_t mask)
+{
+ bc_init_resource(bc, BC_PRIVVMPAGES,
+ BC_PRIVVMPAGES_BARRIER, BC_PRIVVMPAGES_LIMIT);
+ return 0;
+}
+
+struct bc_resource bc_privvm_resource = {
+ .bcr_name = "privvmpages",
+ .bcr_init = bc_privvm_init,
+};
+
+static int __init bc_privvm_init_resource(void)
+{
+ bc_register_resource(BC_PRIVVMPAGES, &bc_privvm_resource);
+ return 0;
+}
+
+__initcall(bc_privvm_init_resource);

```
