

---

Subject: BC: resource beancounters (v6) (with userpages reclamation + configfs)  
Posted by [dev](#) on Thu, 09 Nov 2006 16:42:56 GMT

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

---

MAJOR CHANGES in v6 (see details below):

- configfs interface instead of syscalls (as wanted by CKRM people...)
- added numfiles resource accounting
- added numtasks resource accounting

numfiles and numtasks controllers demonstrate how clean and simple BC interface is.

Patch set is applicable to 2.6.19-rc5-mm1

-----  
Resource BeanCounters (BC).

BC allows to account and control consumption of kernel resources used by \*group\* of processes (users, containers, ...).

Draft BC description on OpenVZ wiki can be found at [http://wiki.openvz.org/UBC\\_parameters](http://wiki.openvz.org/UBC_parameters)

The full BC patch set allows to control:

- kernel memory. All the kernel objects allocatable on user demand and not reclaimable should be accounted and limited for DoS protection.

e.g. page tables, task structs, vmas etc.

- virtual memory pages. BCs allow to limit a container to some amount of memory and introduces 2-level OOM killer taking into account container's consumption.

pages shared between containers are correctly charged as fractions (tunable).

- network buffers. These includes TCP/IP rcv/snd buffers, dgram snd buffers, unix, netlinks and other buffers.

- minor resources accounted/limited by number: tasks, files, flocks, ptys, siginfo, pinned dcache mem, sockets, iptentries (for containers with virtualized networking)

Summary of changes from v5 patch set:

- \* configfs interface instead of syscalls (as wanted by CKRM people)
- \* added numfiles resource accounting
- \* added numtasks resource accounting
- \* introduced dummy\_resource to handle case when no resource registered
- \* calls to rss accounting are integrated to rmap calls

Summary of changes from v4 patch set:

- \* changed set of resources - kmemsize, privvmpages, physpages
- \* added event hooks for resources (init, limit hit etc)
- \* added user pages reclamation (bc\_try\_to\_free\_pages)
- \* removed pages sharing accounting - charge to first user
- \* task now carries only one BC pointer, simplified
- \* make set\_bcid syscall move arbitrary task into BC
- \* resources are not recharged when task moves
- \* each vm\_area\_struct carries a BC pointer

Summary of changes from v3 patch set:

- \* Added basic user pages accounting (lockedpages/privvmpages)
- \* spell in Kconfig
- \* Makefile reworked
- \* EXPORT\_SYMBOL\_GPL
- \* union w/o name in struct page
- \* bc\_task\_charge is void now
- \* adjust minheld/maxheld splitted

Summary of changes from v2 patch set:

- \* introduced atomic\_dec\_and\_lock\_irqsave()
- \* bc\_adjust\_held\_minmax comment
- \* added \_\_must\_check for bc\_\*charge\* funcs
- \* use hash\_long() instead of own one
- \* bc/Kconfig is sourced from init/Kconfig now
- \* introduced bcid\_t type with comment from Alan Cox
- \* check for barrier <= limit in sys\_set\_bclimit()
- \* removed (bc == NULL) checks
- \* replaced memcpy in beancounter\_findcrate with assignment
- \* moved check 'if (mask & BC\_ALLOC)' out of the lock
- \* removed unnecessary memset()

Summary of changes from v1 patch set:

- \* CONFIG\_BEANCOUNTERS is 'n' by default
- \* fixed Kconfig includes in arches
- \* removed hierarchical beancounters to simplify first patchset
- \* removed unused 'private' pointer
- \* removed unused EXPORTS

- \* MAXVALUE redeclared as LONG\_MAX
- \* beancounter\_findcreate clarification
- \* renamed UBC -> BC, ub -> bc etc.
- \* moved BC inheritance into copy\_process
- \* introduced reset\_exec\_bc() with proposed BUG\_ON
- \* removed task\_bc beancounter (not used yet, for numproc)
- \* fixed syscalls for sparc
- \* added sys\_get\_bcstat(): return info that was in /proc
- \* cond\_syscall instead of #ifdefs

Many thanks to Oleg Nesterov, Alan Cox, Matt Helsley and others for patch review and comments.

Thanks,  
Kirill

---