

Containers:

Commodity HW is becoming more powerful. This is giving opportunity to run different workloads on the same platform for better HW resource utilization. To run different workloads efficiently on the same platform, it is critical that we have a notion of limits for each workload in Linux kernel. Current cpuset feature in Linux kernel provides grouping of CPU and memory support to some extent (for NUMA machines).

For example, a user can run a batch job like backup inside containers. This job if run unconstrained could step over most of the memory present in system thus impacting other workloads running on the system at that time. But when the same job is run inside containers then the backup job is run within container limits.

We use the term container to indicate a structure against which we track and charge utilization of system resources like memory, tasks etc for a workload. Containers will allow system admins to customize the underlying platform for different applications based on their performance and HW resource utilization needs. Containers contain enough infrastructure to allow optimal resource utilization without bogging down rest of the kernel. A system admin should be able to create, manage and free containers easily.

At the same time, changes in kernel are minimized so as this support can be easily integrated with mainline kernel.

The user interface for containers is through configs. Appropriate file system privileges are required to do operations on each container. Currently implemented container resources are automatically visible to user space through /configs/container/<container\_name> after a container is created.

Signed-off-by: Rohit Seth <rohitseth@google.com>

Diffstat for the patch set (against linux-2.6.18-rc6-mm2\_:

```
Documentation/containers.txt | 65 ++++
fs/inode.c                   | 3
include/linux/container.h    | 167 ++++++++
include/linux/fs.h           | 5
include/linux/mm_inline.h    | 4
include/linux/mm_types.h     | 4
```

```

include/linux/sched.h      | 6
init/Kconfig               | 8
kernel/Makefile            | 1
kernel/container_configfs.c | 440 ++++++
kernel/exit.c              | 2
kernel/fork.c              | 9
mm/Makefile                | 2
mm/container.c             | 658 ++++++
mm/container_mm.c          | 512 ++++++
mm/filemap.c               | 4
mm/page_alloc.c            | 3
mm/rmap.c                  | 8
mm/swap.c                  | 1
mm/vmscan.c                | 1
20 files changed, 1902 insertions(+), 1 deletion(-)

```

Changes from version 1:

Fixed the Documentation error

Fixed the corruption in container task list

Added the support for showing all the tasks belonging to a container through showtask attribute

moved the Kconfig changes to init directory (from mm)

Fixed the bug of unregistering container subsystem if we are not able to create workqueue

Better support for handling limits for file pages. This now includes support for flushing and invalidating page cache pages.

Minor other changes.

\*\*\*\*\*

This patch set has basic container support that includes:

- Create a container using mkdir command in configfs
- Free a container using rmdir command
- Dynamically adjust memory and task limits for container.
- Add/Remove a task to container (given a pid)
- Files are currently added as part of open from a task that already belongs to a container.
- Keep track of active, anonymous, mapped and pagecache usage of container memory
- Does not allow more than task\_limit number of tasks to be created in the container.

- Over the limit memory handler is called when number of pages (anon + pagecache) exceed the limit. It is also called when number of active pages exceed the page limit. Currently, this memory handler scans the mappings and tasks belonging to container (file and anonymous) and tries to deactivate pages. If the number of page cache pages is also high then it also invalidate mappings. The thought behind this scheme is, it is okay for containers to go over limit as long they run in degraded manner when they are over their limit. Also, if there is any memory pressure then pages belonging to over the limit container(s) become prime candidates for kernel reclaimer. Container mutex is also held during the time this handler is working its way through to prevent any further addition of resources (like tasks or mappings) to this container. Though it also blocks removal of same resources from the container for the same time. It is possible that over the limit page handler takes lot of time if memory pressure on a container is continuously very high. The limits, like how long a task should schedule out when it hits memory limit, is also on the lower side at present (particularly when it is memory hogger). But should be easy to change if need be.

- Indicate the number of times the page limit and task limit is hit
- Indicate the tasks (pids) belonging to container.

Below is a one line description for patches that will follow:

[patch01]: Documentation on how to use containers  
(Documentation/container.txt)

[patch02]: Changes in the generic part of kernel code

[patch03]: Container's interface with configs

[patch04]: Core container support

[patch05]: Over the limit memory handler.

TODO:

- some code(like container\_add\_task) in mm/container.c should go elsewhere.
- Support adding/removing a file name to container through configs
- /proc/pid/container to show the container id (or name)
- More testing for memory controller. Currently it is possible that limits are exceeded. See if a call to reclaim can be easily integrated.
- Kernel memory tracking (based on patches from BC)
- Limit on user locked memory
- Huge memory support

- Stress testing with containers
- One shot view of all containers
- CKRM folks are interested in seeing all processes belonging to a container. Add the attribute show\_tasks to container.
- Add logic so that the sum of limits are not exceeding appropriate system requirements.
- Extend it with other controllers (CPU and Disk I/O)
- Add flags bits for supporting different actions (like in some cases provide a hard memory limit and in some cases it could be soft).
- Capability to kill processes for the extreme cases.
- ...

This is based on lot of discussions over last month or so. I hope this patch set is something that we can agree and more support can be added on top of this. Please provide feedback and add other extensions that are useful in the TODO list.

Thanks,  
-rohit

-

To unsubscribe from this list: send the line "unsubscribe linux-kernel" in the body of a message to majordomo@vger.kernel.org  
More majordomo info at <http://vger.kernel.org/majordomo-info.html>  
Please read the FAQ at <http://www.tux.org/lkml/>

---