# CFQ vs Containers

2008/2/26

NTT Open Source Software Center
Fernando Luis Vázquez Cao
Hiroaki Nakano

1. Block I/O resources and cgroups
2. Cgroups
3. I/O group scheduling
4. ioband

# 1．Block I/O resources and cgroups

**State of things**

> ➢ CFQ's IO priority is an attribute of a process so it affects all devices it sends I/O requests to

> ➢ I/O priority can be set by PID, PGRP, or UID, but...

> ➢ ...all the processes that fall within the same class/priority are scheduled together

- Goals

  - Being able to define arbitrary groupings of processes and...

  - ...treat each group as a single scheduling entity

  - Provide (soft) data rate guarantees

  - Perform I/O bandwidth control independently on each device

  - Scheduler-independent I/O bandwidth control

  - Usable even when the generic make_request_fn function is not used

- **What kind of things can be done?**

  - I/O prioritization
    - ionice-like approach

  - Proportional bandwidth scheduling
    - Each process/group of processes has a weight that determines the share of bandwidth they receive

  - I/O limiting
    - Set an upper limit to the bandwidth a group of tasks can use
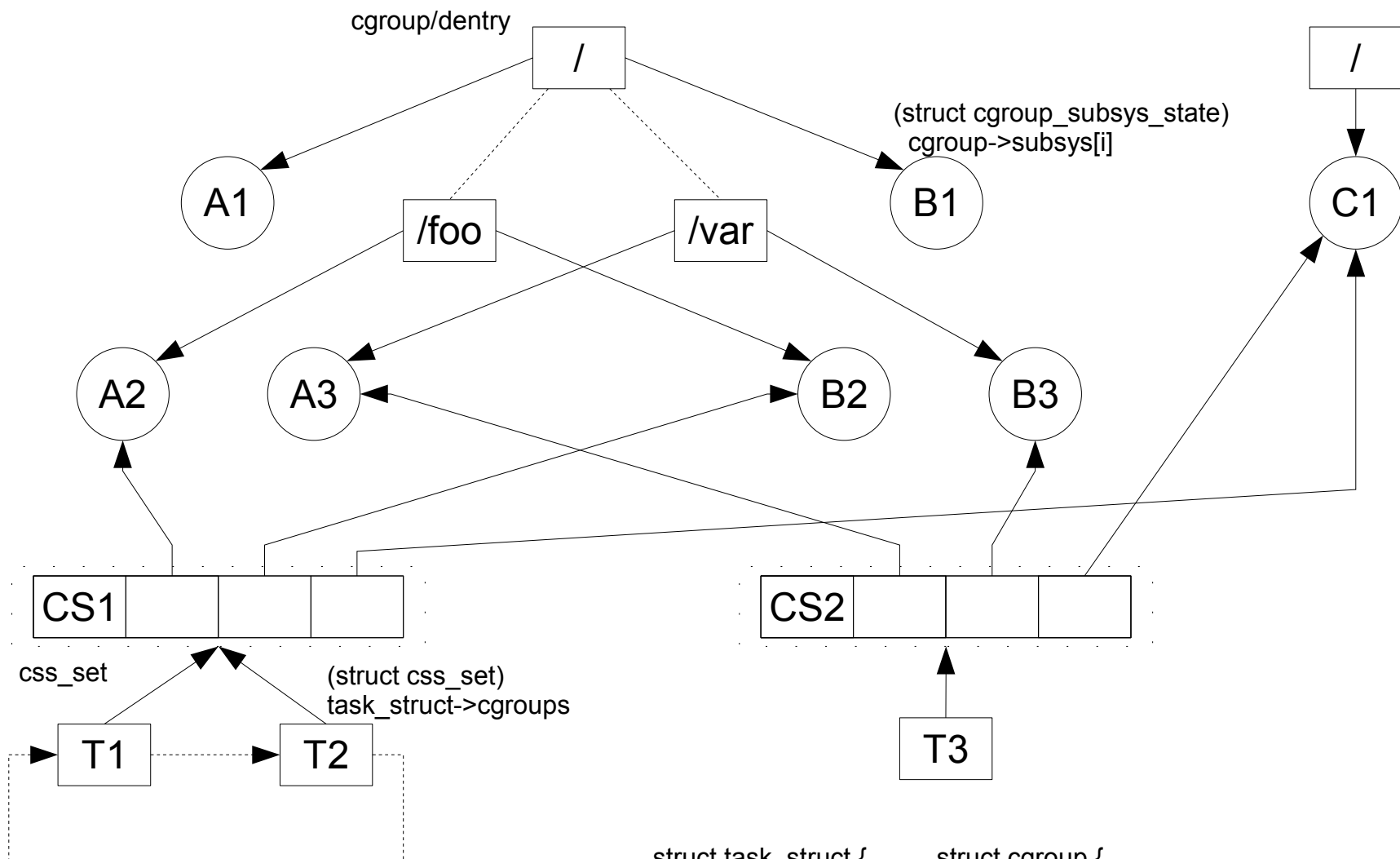
# 2．Cgroups

## Subsystem/controller

➢ Is a part of the kernel, commonly a system resource, which might have an interest in what a group of processes are doing

## Cgroup

➢ Is a group of processes that share a set of parameters used by one or more subsystems

➢ Characteristics

  ✗ Cgroups are hierarchical

  ✗ Each cgroup hierarchy is controlled through a cgroup filesystem whose tree of directories follows the structure of the cgroup hierarchy

cgroup/dentry

/

(struct cgroup_subsys_state)
cgroup->subsys[i]

/

A1

/foo

/var

B1

C1

A2

A3

B2

B3

CS1

CS2

css_set

(struct css_set)
task_struct->cgroups

T1

T2

T3

cg_list (anchored at css_set->tasks)

```
struct task_struct {        struct cgroup {
  ...                         struct list_head sibling;
  css_set *cgroups;           struct list_head children;
  ...                         struct cgroup *parent;
}                             struct dentry *dentry;
                              struct cgroup_subsys_state *subsys[];
                              struct list_head css_sets;
                            }
```
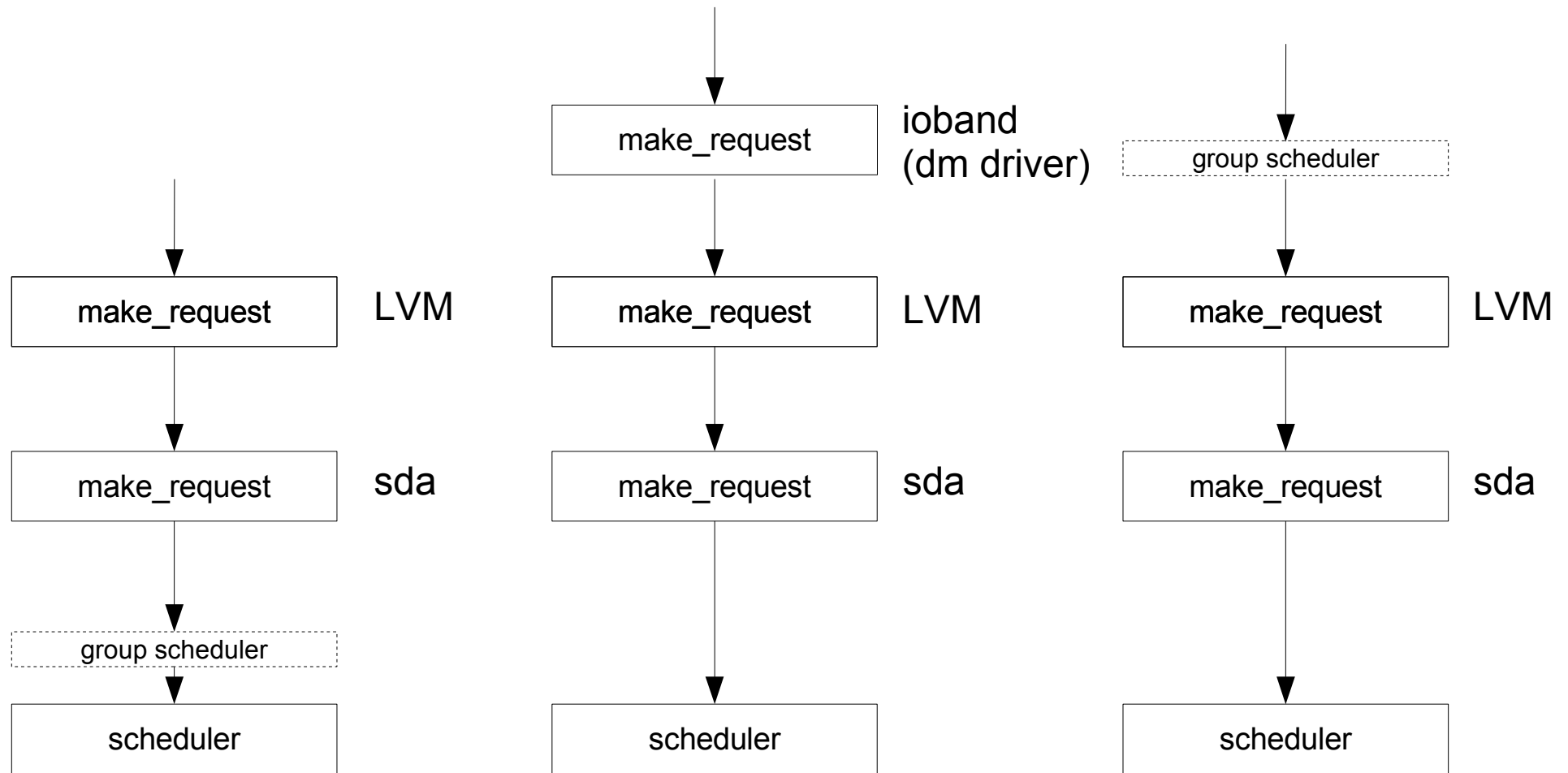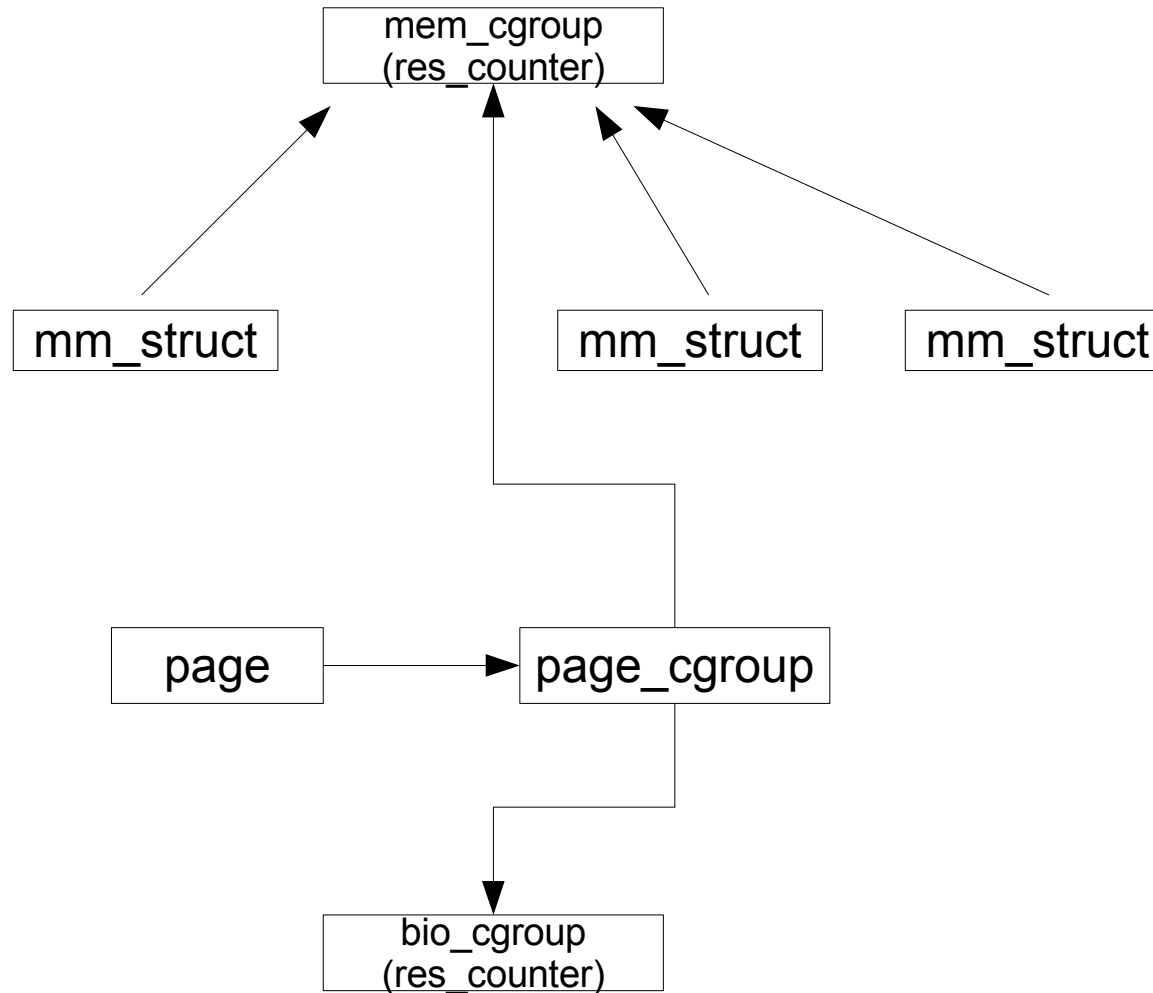
*Original figure by Paul Menage*

# 3・I/O group scheduling

- Cgroups-aware I/O scheduling
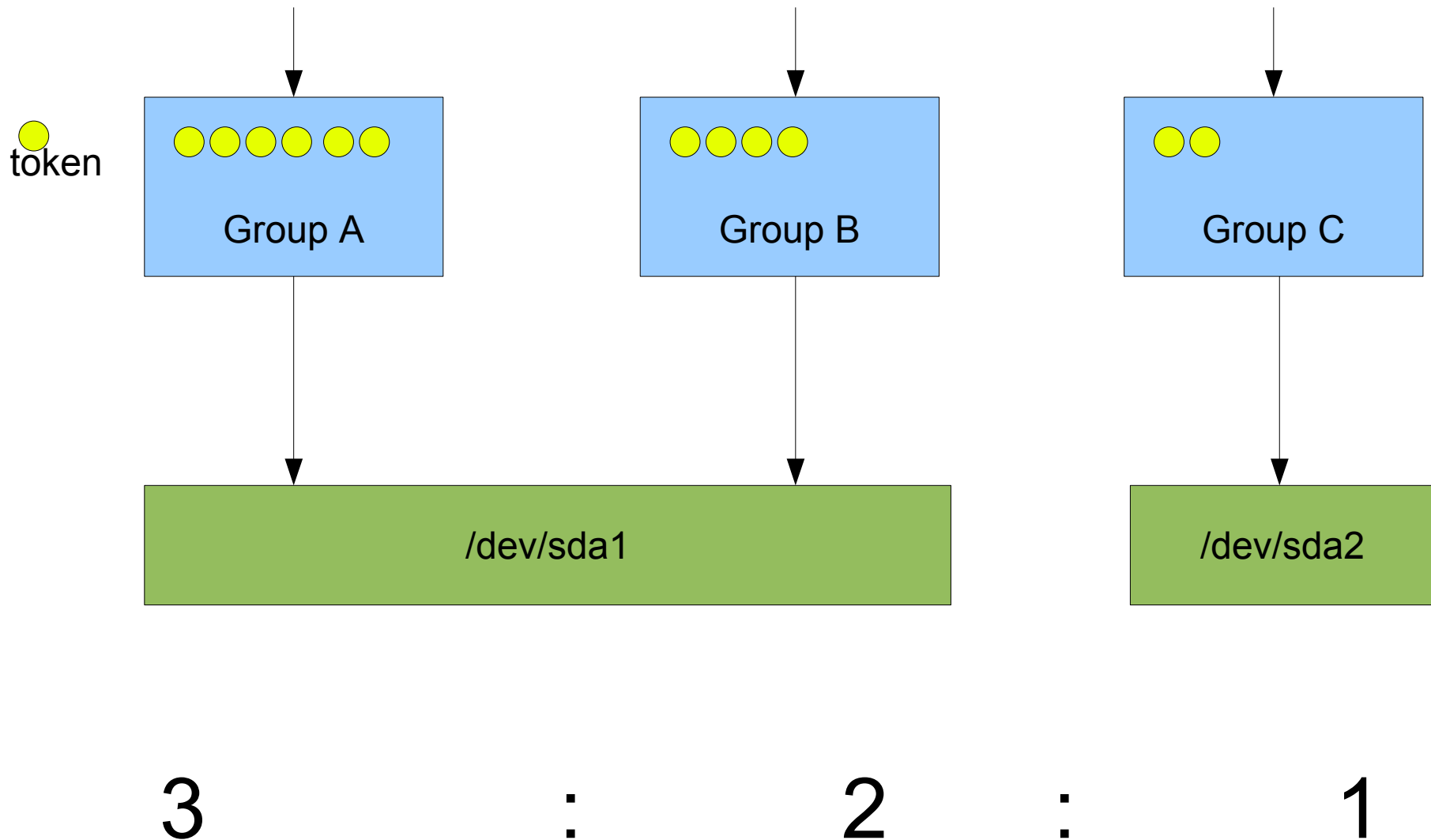
- I/O tracking

- Group scheduling algorithm

token

Group A

Group B

Group C

/dev/sda1

/dev/sda2

3 : 2 : 1

# 4 . ioband

✗ I/O bandwidth controller implemented as a device-mapper driver
✗ Bandwidth assigned according to the relative weight of each job

| cgroup A | cgroup B | the others |
|---|---|---|

jobs

| group | group | default group |
|---|---|---|

ioband groups

ioband 1

ioband devices

| PID X | PID Y | the others |
|---|---|---|

| group | group | default group |
|---|---|---|

ioband 2

| sdb1 | sdb2 |
|---|---|

physical devices

- **Pros**
  - ➢ Works with any I/O scheduler
    - ✗ This is a direct consequence of using a dm driver for the implementation
  - ➢ Each device can be configured independently
    - ✗ As opposed to CFQ's I/O priority which affects all I/O generated by a process

- **Cons**
  - ➢ Overkill?

# Thanks for your attention

Contact:   fernando@oss.ntt.co.jp

nakano.hiroaki@oss.ntt.co.jp