# News and trends in Linux 2.6

Thomas Petazzoni
Free Electrons
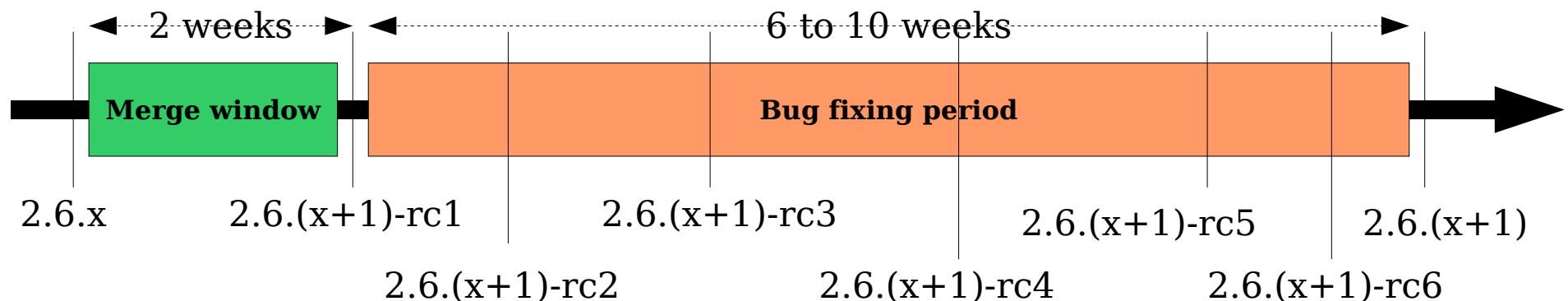http://free-electrons.com/

**Free Electrons**

# Development model

- The previous development/stable versions development model has been dropped

- The 2.6 branch exists since 2003, and no 2.7 development branch has been opened

- New features are gradually added in each version of the kernel

  - Fits better the distributors needs (no need to backport a lot of features from development kernels)

  - Nicer for kernel developers : their changes are sooner made available to users

**Free Electrons**

# Development model (2)

► After the release of a 2.6.x kernel

  ► Two weeks of *merge window* to integrate new features, drivers, etc. Closed by the release of 2.6.(x+1)-rc1

  ► Six to ten weeks bug-fixing period, with regular releases of 2.6.(x+1)-rcY kernels

► Release of 2.6.(x+1) once considered stable

```
     2 weeks              6 to 10 weeks

    Merge window                  Bug fixing period

2.6.x        2.6.(x+1)-rc1        2.6.(x+1)-rc3        2.6.(x+1)-rc5        2.6.(x+1)
                2.6.(x+1)-rc2        2.6.(x+1)-rc4        2.6.(x+1)-rc6
```

**Free Electrons**

3

# Source code management tool (1)

▶ Until 2002 : no central source code management tool

▶ In 2002, switch to BitKeeper

    ▶ During the 2.5 development cycle

    ▶ BitKeeper was a distributed SCM, that suited well the development model of the kernel. Free version available, but proprietary.

▶ In April 2005, announcement that BitKeeper will no longer be free

    ▶ Torvalds starts the development of a new distributed SCM, Git

    ▶ Git is now used by the kernel and many other free software projects

    ▶ Maintained by Junio Hamano.

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com
Jul 8, 2008

**Free Electrons**

4

# Source code management tool (2)

- Git is now used by most core kernel developers

- Allows to get an higher number of patches integrated

  - Kernel development is going on with a very fast rhythm

  - 5000 to 7000 patches in each kernel version

  - 85 lines added to the kernel every hour, since 3 years

- Distributed nature well-suited for free software projects

  - Local branches for development

  - Exchange development versions between developers

  - Offline development
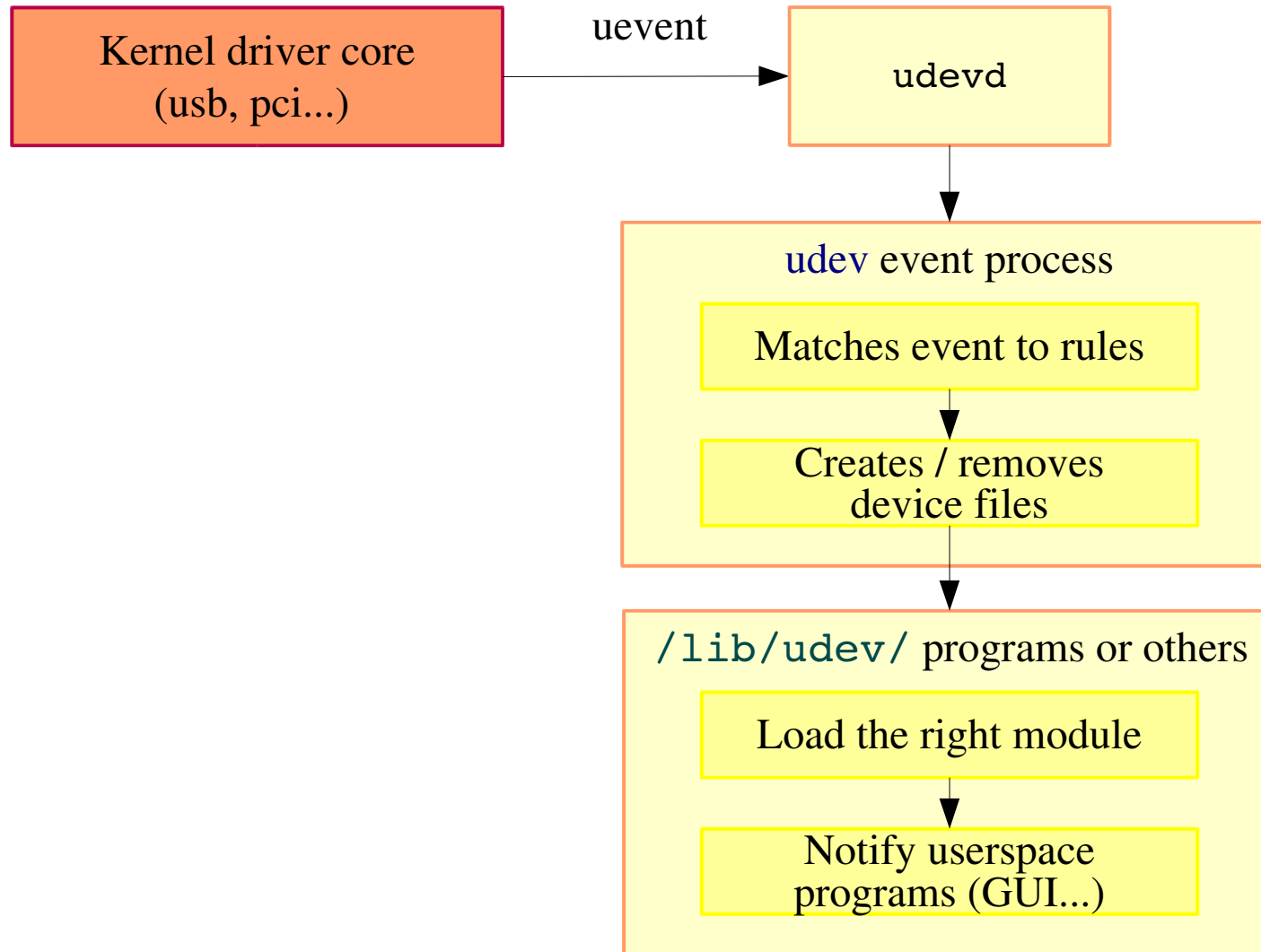
  - Powerful merging facilities

**Free Electrons**

# Kobject events userspace notification

- Allow userspace to be notified of kernel events occurring on kobjectfs

- Uses a netlink socket

- Userspace can be notified of events such as

  - Device insertion and removal

  - Mount notifications

  - Simple events (CPU overheating, etc.)

- Foundation for udev

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com
Jul 8, 2008

**Free Electrons**

6

# udev - typical operation

```
Kernel driver core          uevent
(usb, pci...)        ─────────────────▶         udevd
```

**udev event process**

```
Matches event to rules
        │
        ▼
Creates / removes
   device files
```

**/lib/udev/ programs or others**

```
Load the right module
        │
        ▼
Notify userspace
programs (GUI...)
```

*Free Electrons*

Jul 8, 2008

# devfs removal

▶ Devfs was a virtual filesystem added in the 2.5 development branch

▶ Its goal was to provide a dynamic `/dev` directory, instead of having thousands of static and useless entries for every possible device

▶ For several reasons, it was disliked by many kernel developers

    ▶ Broken naming

    ▶ Naming policy in the kernel

    ▶ Not flexible enough

▶ Finally removed in 2.6.18, superseded by udev.

**Free Electrons**

# CPU sets

- Big NUMA systems present challenges for the efficient scheduling and memory placement of processes

- On small systems, CPU affinity might be sufficient, but not on big NUMA systems

- Allows to create cpusets, consisting of

  - CPU nodes

  - Memory nodes

- And then to assign tasks to a given cpuset

- User interface in the form of a virtual cpuset filesystem

**Free Electrons**

# inotify

- Provide filesystem events notification to userspace programs

- Replacement for `dnotify`

  - Relied on one file descriptor for each monitored directory. Issues with the maximum number of file descriptors, and removable devices, awful userspace interface

  - Can only monitor changes at the directory level, not file level

- Notification of events such as access, modification, change of attributes, open, close, move, rename, creation, deletion

- Mostly used for desktop search system so that reindexing is not needed

- Might be useful for other applications as well

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com
Jul 8, 2008

**Free Electrons**

10

# FUSE

- Filesystem in userspace

- Allows to write filesystem drivers in userspace, using `libfuse`

- These filesystems can be used like traditional filesystem: mounted and accessed using the standard API

- The fuse kernel module *forwards* VFS calls to an userspace program

- Wide variety of userspace filesystem drivers: sshfs, ntfs-3g, zfs, gmailfs, etc.

**Free Electrons**

Jul 8, 2008

# New system calls (1)

- ► `*at()` system calls

  - ► `openat(), mkdirat(), mknodat(), fchownat(), futimesat(), fstatat(), unlinkat(), renameat(), linkat(), symlink_at(), readlinkat(), fchmodat(), faccessat()`

  - ► Operations relative to an open file descriptor corresponding to a directory

  - ► Needed to implement race-free filesystem traversal

  - ► Needed to implement virtual per-thread current directory, for backup software.

**Embedded Linux System development**

Free Electrons

© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com
Jul 8, 2008

12

# New system calls (2)

- ▶ `pselect(), ppoll()`

    - ▶ Signal-aware version of `select()` and `poll()`

    - ▶ No more race condition or complicated self-pipe trick needed to properly handle signals with `select()` and `poll()`

    - ▶ See http://lwn.net/Articles/176911/

- ▶ `unshare()`

    - ▶ Allows to selectively unshare resources that are normally shared between threads of the same process

    - ▶ Polyinstantiated directories, for security

    - ▶ Intra-application isolation

    - ▶ See Documentation/unshare.txt

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com

*Jul 8, 2008*

**Free Electrons**

**13**

# New system calls (3)

- `tee()`, `splice()`, `vmsplice()`

  - `splice()` allows to moves data from one file descriptor to another file descriptor without involving userspace

  - `tee()` does the same, except that it doesn't consume the input data, it can still be read from userspace

  - `vmsplice()` allows to feed userspace memory into a pipe

  - Increase performances by decreasing the number of memory copies

  - See `tee(2)`, `splice(2)` and `vmsplice(2)`

- `sync_file_range()`

  - Allows to synchronize part of a file to the disk, and wait for the completion of part of the process, depending on given flags.

**Free Electrons**

▶ `fallocate()`

> ▶ Allows to preallocate or deallocate blocks on disk
>
> ▶ Ensure that blocks are contiguous on disk, which can improve performances
>
> ▶ Only implemented in ext4 and OCFS2 so far

**Free Electrons**

Jul 8, 2008

**15**

# SLOB / SLUB allocators

- Historically, the Linux kernel relies on an allocator using the SLAB strategy for small allocations

- In 2.6, two new alternative allocators, also implementing the SLAB strategy have been added

  - SLOB, with a focus on code size and low memory overhead, for embedded systems

  - SLUB, with a focus on scalability while remaining general-purpose, now the default allocator

**Free Electrons**

# configfs

- New virtual filesystem

- Aimed at easing the configuration of complex kernel applications from userspace

- When directories and files are created by userspace applications, callbacks are called in the kernel

- Allows greater flexibility of configuration, without using ioctl() calls

- Used by DLM, OCFS and netconsole for the moment

- See http://lwn.net/Articles/148973/

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com
Jul 8, 2008

**Free Electrons**

**17**

# Clustering

- OCFS 2

  - Clustering filesystem developed by Oracle. Allows several hosts to access the same storage by taking care of concurrent accesses

- GFS

  - Clustering filesystem originally developed by Sistina, bought by RedHat and released under the GPL

- TIPC

  - Transparent Inter Process Communication protocol

  - Intra cluster communication

  - Originates from Ericsson

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com

*Free Electrons*

Jul 8, 2008

**18**

# High resolution timers

- Classic timers in Linux have at most a resolution equal to the timer frequency (1 ms, 4 ms or 10 ms)

  - Resolution not sufficient for multimedia or realtime applications

- High resolution timers allows to have a much higher resolution, depending on what the hardware allows

- Provide a new in-kernel API, also used to improve the userspace time-related functions

  - Nanosleep

  - Itimers

  - POSIX timers

# Userspace priority inheritance

- Priority inversion occurs when a low-priority process holds a lock needed by an higher-priority process

- A classical solution is to temporarily boost the low-priority process's priority to the high-priority process's priority, until it releases the lock : priority inheritance

- Support for priority inheritance in userspace locks has been added in 2.6.18, through the *futex* facility

- Can be used using
  ```
  pthread_mutexattr_setprotocol (...,
  PTHREAD_PRIO_INHERIT)
  ```

**Free Electrons**

# ext4 (1)

- New major developments had to be made on the ext3 filesystem

  - Did not match the stability requirements of ext3

  - A new ext4 filesystem has been created, initially a simple copy of ext3

- Goal is to improve the scalability and reliability of the ext filesystem

  - Large filesystems (64 bits)

  - Challenge of the growing size of hard disks but not their throughput or access time

Jul 8, 2008

**Free Electrons**

- Features

  - Extents

  - Filesystems bigger than 16 TB

  - Internal redundancy to improve behavior in the face of corruption

  - Improved file allocation

  - Large inodes (nanoseconds timestamps, etc.)

  - Reduced fsck time

- Not production-ready

- See Documentation/filesystems/ext4.txt

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com

**Free Electrons**

Jul 8, 2008

**22**

# Namespaces

- One method of virtualization is isolation

    - Single kernel, but processes are enclosed in isolated environments

    - Approach used by vserver and OpenVZ

- Needs support from all the kernel subsystems

    - User namespaces, 2.6.23

    - PID and network namespaces, 2.6.24

    - New `clone()` flags

    - More subsystem namespaces in next versions of the kernel

- Also useful for checkpoint and restart of processes in HPC

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com

**Free Electrons**

Jul 8, 2008

23

# Paravirtualization

- Paravirtualization

  - Running a guest kernel under a host kernel, where the guest has been modified to run on a virtual architecture, similar to the hardware architecture

  - The guest kernel must call an hypervisor for privileged operations (enable/disabling interrupts, for example)

- Generic paravirtualization support, introduced in 2.6.20

  - Encapsulate the set of operations for which a call to an hypervisor has to be made into a set of functions pointers, *paravirt_ops*

  - Defaults to their normal behavior, but can be modified by the hypervisor

- Generic paravirtualization support, 2.6.20

- Xen, 2.6.23

- Lguest, 2.6.23

# Paravirtualization (2)

- Xen, 2.6.23

  - Parts of Xen have been merged

  - Only guest support, no hypervisor, no dom0

  - Based on `paravirt_ops`

- Lguest, 2.6.23

  - Simple hypervisor for Linux on Linux

  - Demonstrate the how powerful the `paravirt_ops` infrastructure is

  - Hypervisor in a kernel module

  - Specifically-compiled guest kernels

  - See Documentation/lguest

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com
Jul 8, 2008

**Free Electrons**

25

# Paravirtualization (3)

- Virtual Machine Interface, 2.6.21

  - Proposed by VMware

  - The `paravirt_ops` approach was preferred

  - The port of VMI on top of `paravirt_ops` has been merged in 2.6.21

  - Another interface between guest kernels and an hypervisor

  - Allows Linux guest kernels to run on VMware products

# KVM

▶ Driver for Intel and AMD hardware virtualization extensions

▶ Export them to userspace through `/dev/kvm`

▶ Using this driver, an userspace process can run a guest virtual machine with fully virtualized hardware

▶ Virtualizes CPU and memory,
but the rest of the hardware has to be emulated

▶ Works in combination with Qemu

▶ Runs unmodified guests

▶ Also supports a paravirtualized mode, migration, etc.

▶ See http://kvm.qumranet.com/

**Embedded Linux System development**

© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com

*Free Electrons*

Jul 8, 2008

27

# Wireless stack

- Huge differences between wireless adapters in the range of functionalities supported in hardware and those that have to be implemented in software

  - Mess and duplication in Linux drivers

- The Devicescape company developed and released under the GPL a new wireless stack, now integrated to the kernel

- Complete software MAC, WEP, WPA, link-layer bridge module, hostapd, QoS, etc.

- New kernel/user interface, based on netlink instead of ioctl()

- Kernel drivers are being rewritten, but will be cleaner and simpler thanks to the new stack.

*Free Electrons*

# Firewire stack

- New Firewire stack being developed

- *« get a small, maintainable and supportable FireWire stack »*

- Smaller code base and binaries

- Cleaned-up and improved in-stack APIs and design

- Consolidation of the four userspace ABIs into a single improved ABI

- Solve bugs and security issues

- Compatibility at the library level (libraw1394, libdc1394)

- Both stacks are in the kernel, until the new one fully replaces the old one.

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com
Jul 8, 2008

**Free Electrons**

**29**

# Libata PATA

- New SATA devices and controllers are very similar to SCSI devices and controllers

  - From the origin, SATA support has been integrated inside the SCSI stack, using libata

- libata PATA is an effort to propagate this idea to traditional PATA devices and controllers

- In the end, all devices should appear as SCSI devices, with a single API

- Production-ready, already used on modern distributions

- See Documentation/DocBook/libata

**Free Electrons**

# OSS removal

- Open Sound System

    - Original sound subsystem of the Linux kernel

    - Now superseded by ALSA

- OSS drivers are slowly removed, when an ALSA driver exists for the same driver

- Should be completely dropped in a couple of years

- New drivers should be developed inside the ALSA framework

# UBI

- « Unsorted Block Images »

- Volume management system for flash devices

    - Manages multiple volumes on a single flash

    - Spreads the I/O load across the whole flash chip (global wear-leveling)

    - Transparently handles bad physical erase blocks

    - Handles bit-flips in the Flash by remapping blocks when bit-flips are detected

    - Atomic logical eraseblock change

- See http://www.linux-mtd.infradead.org/doc/ubi.html

# Tick frequency

▶ In 2.4, the tick frequency was 100 Hz

▶ In 2.6, it was increased to 1000 Hz

▶ In 2.6.13, it was reduced to 250 Hz, but made configurable

▶ In 2.6.21, the Dynamic Ticks patch is merged

   ▶ It removes the regular timer tick

   ▶ Relies on programming the hardware timer so that the system is only
   woken up for the next event

   ▶ Allows to reduce power consumption

   ▶ Initial support only for x86

   ▶ Extended to x86_64, PPC, ARM and MIPS in 2.6.24

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com

*Free Electrons*

Jul 8, 2008

33

# CFS (1)

- During the 2.5 cycle, an O(1) scheduler, developed by Ingo Molnar, was merged

  - Scalability improvements with many processes

- Issue with O(1) scheduler and interactive tasks, that did not react quickly enough in presence of CPU-bound tasks

  - Interactivity estimator and heuristics developed by Con Kolivas

- Con Kolivas developed a brand new scheduler, called RDSL, based on a strict fairness

- Ingo Molnar took over the idea, and developed CFS, Completely Fair Scheduler, which was finally merged.

# CFS (2)

- In the O(1) scheduler, two arrays of run queues to keep track of processes

  - Each array with 140 entries to take priorities into account

- In the CFS, all runnable processes are stored in a red-black tree

- The key in the red-black tree is the time
that the process should have had on the CPU

  - Basically the time for which the process waited without being run, divided by the number of processes, with a correction for priority

  - Allows to be perfectly fair with processes, and be nice with interactive processes

- Introduction of modular, chained, schedulers

**Free Electrons**

# Lumpy reclaim, anti-fragmentation

► Fragmentation of physical memory quickly becomes an issue for big allocations on highly-loaded machines

► Several improvements made in this area

  ► Lumpy reclaim

    ► Tune the page reclaim algorithm to not only take the recency into account, but also the fact that pages are grouped or not

  ► Anti-fragmentation

    ► Grouping pages of related type together (movable, reclaimable, unreclaimable, highatomic)

    ► When pages are migrated or reclaimed, large regions of physical memory are freed, allowing bigger allocations to succeed

Jul 8, 2008

**Free Electrons**

# UIO

- Framework that allows the development of device drivers in userspace

- Small kernel module to register the device and for the interrupt handler, the rest in userspace

- Easier to debug

- Driver can be kept proprietary with no licensing issues

- Interface: `/dev/uioX`

  - `mmap()` to read/write registers

  - Blocking `read()` to be notified of interrupts

- http://www.free-electrons.com/kerneldoc/latest/DocBook/uio-howto/

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com
*Jul 8, 2008*

*Free Electrons*

37

# Power management

▶ **Suspend and resume**

▶ **Cpuidle**, generic framework for supporting software-controlled
   idle processor power management

   ▶ Hardware specific drivers, various governors for the state transition decisions

▶ **Cpufreq**, generic framework for frequency and voltage scaling

   ▶ Hardware specific drivers, architecture-independent governors

▶ **PmQOS**

   ▶ Framework for expressing latency constraints, and make sure that they are
      taken into account for power management decisions

**Free Electrons**

# Linux Kernel Markers

- Static probing points inserted at various points inside the kernel

- Well-defined trace points at correct places in the kernel

  - Added by subsystem developers

  - Maintained in the official kernel tree

- Allow users and system administrators to use higher-level tools

  - SystemTap

  - LTTng

- See http://lwn.net/Articles/245671/ and Documentation/markers.txt

Jul 8, 2008

# Fault injection

- Framework for injecting faults in a running kernel

- The goal is to test error code paths

- Allows to

    - Make SLAB allocations fail

    - Make page allocations fail

    - Return I/O errors

- Userspace configuration through debugfs

    - Set the rate of faults

    - Filter the injection to certain tasks or particular pieces of code

- See Documentation/fault-injection/fault-injection.txt

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com

*Jul 8, 2008*

**Free Electrons**

**40**

# blktrace

- Kernel module and userspace tool to trace I/O requests

- All the steps of an I/O request can be traced

- Useful to debug issues with I/O loads

- Relies on relayfs to transfer large amount of informations from the kernel to userspace

```
% blktrace -d /dev/sda -o - | blkparse -i -
  8,0    3        1     0.000000000   697 G W 223490 + 8 [kjournald]
  8,0    3        2     0.000001829   697 P R [kjournald]
  8,0    3        3     0.000002197   697 Q W 223490 + 8 [kjournald]
  8,0    3        4     0.000005533   697 M W 223498 + 8 [kjournald]
  8,0    3        5     0.000008607   697 M W 223506 + 8 [kjournald]
  8,0    3        6     0.000011569   697 M W 223514 + 8 [kjournald]
```

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com

*Free Electrons*

Jul 8, 2008

41

# Lockdep

▶ Kernel lock validator

▶ Allows to detect deadlock before they occur, even rare deadlocks

▶ Associate each spinlock with a key, so that similar locks are handled only once

▶ When locking, look at all already taken locks, and make sure that none of these locks are taken after the newly taken lock in other contexts.

▶ When unlocking, make sure that the lock being unlocked is at the top of the taken locks.

▶ Validate spinlocks vs interrupts behavior.

# PowerTOP

▶ With dynamic tick, new interest in fixing the kernel code and applications that wake up the system too often

▶ PowerTOP application allows to track the worst offender

**Embedded Linux System development**

© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com

*Jul 8, 2008*

**Free Electrons**

**43**

# LatencyTop support

- Userspace tool and kernel support to detect where the latencies are

- Accumulate the time spent sleeping between two points, to find where the latencies come from

- See http://www.latencytop.org/

```
jake@bike:/h/jake/latencytop-0.3 - Shell - Konsole

Session   Edit   View   Bookmarks   Settings   Help

   LatencyTOP version 0.3       (C) 2008 Intel Corporation

Cause                                      Maximum        Average
Unlinking file ()                          82.6 msec      1.5 msec
Synchronous bufferhead read ()             55.3 msec      3.4 msec
EXT3 reading inode ()                      40.6 msec      3.5 msec
Reading from file ()                       21.5 msec      4.9 msec
process fork ()                            18.8 msec      0.2 msec
SCSI layer command execute ()              14.6 msec      1.2 msec
page fault ()                              14.2 msec      2.0 msec
Unknown reason (do_wait+0xa0b/0xb15)       13.2 msec      0.8 msec
Writing to file ()                         10.0 msec      5.1 msec




Process hald-addon-stor (3242)
SCSI layer command execute ()              14.6 msec      1.2 msec
SCSI cdrom ioctl ()                         8.1 msec      0.8 msec








 scsi_eh_0   kjournald   kjournald   kjournald   hald-addon-inpu  hald-addon-stor

   Shell
```

**Free Electrons**

# New architectures

- FRV

- AVR32

- Blackfin

- Xtensa

- MN10300/AM33

**Free Electrons**

# New Drivers

ATI r500 DRM, rs690 DRM, Sis 662/671, BF52x EZkit Display, S3c2410 fb, Palmchip BK3710 IDE, Sega Dreamcast GD-ROM, RB500 PATA Compactflash, Marvell 6121 SATA, HT1100 SATA, Onkyo SE-90PCI and SE-200PCI, Xilinx ML403 AC97, ASoC TLV320AIC3X, SiS 7019 Audio Accelerator, Kore controller 2, Asus Xonar, CMI8788, ASoC TLV, S3c2412 IIS, ASoC drivers for the Freescale MPC8610 SoC, Marvell 6440 SAS/SATA, enc28j60, ath5k, RDC R6040 Fast Ethernet, bnx2x for BCM57710, tuner-xc2028, MT9V111, DViCO FusionHDTV Dual Digital 4, DViCO FusionHDTV NANO2 w/ZL10353, AVerMedia EZMaker PCI Deluxe, cs5345, Xceive xc5000 silicon tuner, Hauppauge HVR1500Q, NXP TDA18271HD/C, Add Beholder TV 401/405/407/409/505/507/609/M6, Syntek DC1125 webcams, Pinnacle 800i, Creative DiVi CAM 516, Twinhan Hybrid DTV-DVB 3056 PCI, Medion / Creatix CTX948, Genius TVGo A11MCE, VT8237S, PCF8575 chip, Apple aluminum USB keyboards, Genius KB-29E, Logitech Elite keyboards, Winchiphead USB->RS 232, Aircard 881U, Onda H600/Zte MF330, RATOC REX-USB60F, iuu_phoenix driver, SH7722 USBF, ohci-sm501 driver, Motorola ROKR Z6, Neteffect RNICs, Texas Instruments/Burr-Brown ADS7828, Winbond W83L786NG/NR, Analog Devices ADM1024, Analog Devices ADT7473, pasemi_nand driver, SST 39VF1601... (incomplete)

## Just in 2.6.25 !

*Free Electrons*

Jul 8, 2008

46

# Soon in a kernel near you...

**Embedded Linux System development**

© Copyright 2004-2008, Free Electrons

Creative Commons Attribution-ShareAlike 2.5 license

http://free-electrons.com

Jul 8, 2008

**Free Electrons**

**47**

# KGDB

http://kgdb.linsyssoft.com/

▶ The execution of the kernel is fully controlled by `gdb` from another machine, connected through a serial line.

▶ Can do almost everything, including inserting breakpoints in interrupt handlers.

▶ A simplest version of Linsyssoft patch is being worked on by Ingo Molnar for integration inside the kernel, but Linus Torvalds is known not to like debuggers a lot.

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com

*Free Electrons*

*Jul 8, 2008*

**48**

# Filesystems (1)

- BTRFS

  - New generation filesystem developed by Oracle to compete with Sun's ZFS (checksumming, snapshot, volumes, mirroring, striping, extent-based, etc.)

  - Can be associated with CRFS to export filesystems over the network

  - See http://oss.oracle.com/projects/btrfs/

- LogFS

  - Flash filesystem, aimed at replacing JFFS2

  - Focus on scalability, usable on large devices

  - See http://www.logfs.org/logfs/

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com

**Free Electrons**

Jul 8, 2008

49

# Filesystems (2)

- UnionFS

  - Stackable unification filesystem

  - Merge the contents of several directories, while keeping their physical contents separate

  - Useful to merge a small read-write device over a read-only device, for example for LiveCDs or embedded devices

  - See http://www.am-utils.org/project-unionfs.html

- ChunkFS

  - Experimental, chunk-based filesystem, designed to work properly on vast amounts of storage

  - See http://www.valhenson.org/chunkfs/

# Utrace

- Infrastructure for tracing and controlling user threads

- Foundation for writing tracing engines, which can be loaded as modules

- Provides three facilities

  - Thread event reporting

  - Core thread control

  - Thread machine state access

- Replacement for ptrace(), with a similar userspace API

**Free Electrons**

Jul 8, 2008

# Real Time improvements

- Real-time mutex code
    - Turn all spinlocks to mutexes to make them preemptible
- Threaded interrupt handlers
    - Interrupt handlers running in threads allows them to be scheduled and prioritized like any other threads
    - Allows low-priority interrupts to not cause latencies on high-priority tasks
- Other latency reduction patches
- More than 400 patches still in -rt

# Syslets

▶ Syslets are a mechanism to perform asynchronous system calls from userspace

▶ Instead of developing an ad-hoc asynchronous version of each system call, syslets is a generic solution

▶ Executes a system call, and if it happens to block, creates a new thread to continue execution in userspace while the system call is being executed in the original thread

▶ Complex version proposed initially, with multiple syscalls in one functionality.

▶ Simpler version being proposed now, to ease inclusion

**Embedded Linux System development**
© Copyright 2004-2008, Free Electrons
Creative Commons Attribution-ShareAlike 2.5 license
http://free-electrons.com

**Free Electrons**

Jul 8, 2008

**53**

# The big trends

- Scalability

  - In CPU and memory

  - In storage

- Virtualization

- Real-time

- Power management

- Debugging tools

- ... and of course more architectures, more platforms, more drivers

**Free Electrons**

Jul 8, 2008

# Follow the news and trends

- **Linux Weekly News**, http://lwn.net

  - Weekly-published newsletter about free software

  - Excellent coverage of kernel stuff, written by Jonathan Corbet

  - $5 monthly subscription for fresh news, all articles available for free after one week

- **LinuxChanges**, http://kernelnewbies.org/LinuxChanges

  - The human readable kernel ChangeLog

**Free Electrons**

# Questions ?

**Embedded Linux System development**

© Copyright 2004-2008, Free Electrons

Creative Commons Attribution-ShareAlike 2.5 license

http://free-electrons.com

Jul 8, 2008

**Free Electrons**

**56**