

Chasing the Penguin: State and Evolution of the Kernel

Wolfgang Mauerer

MPRG IOIP & linux-kernel.net

10. October 2008

Dynamics of Kernel Development

Kernel Documentation

Understanding the Kernel

- Documenting new Features

- Analysis Tools

Social Aspects

3 Outline

Dynamics of Kernel Development

Kernel Documentation

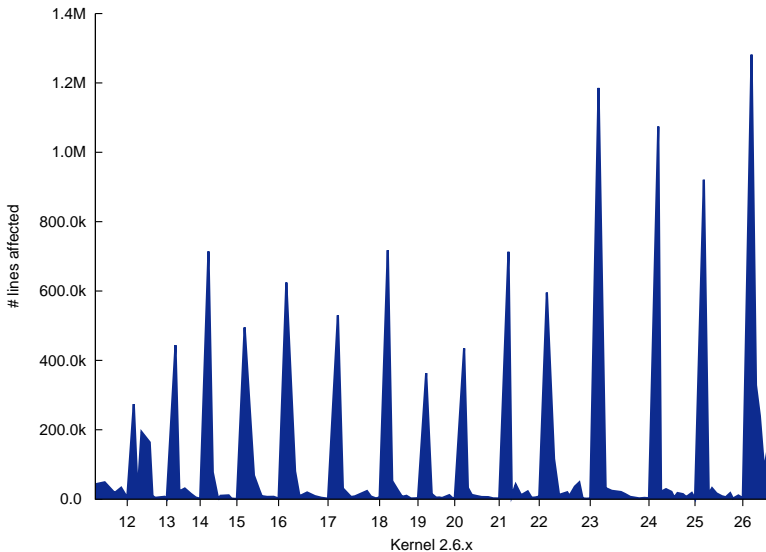
Understanding the Kernel

- Documenting new Features

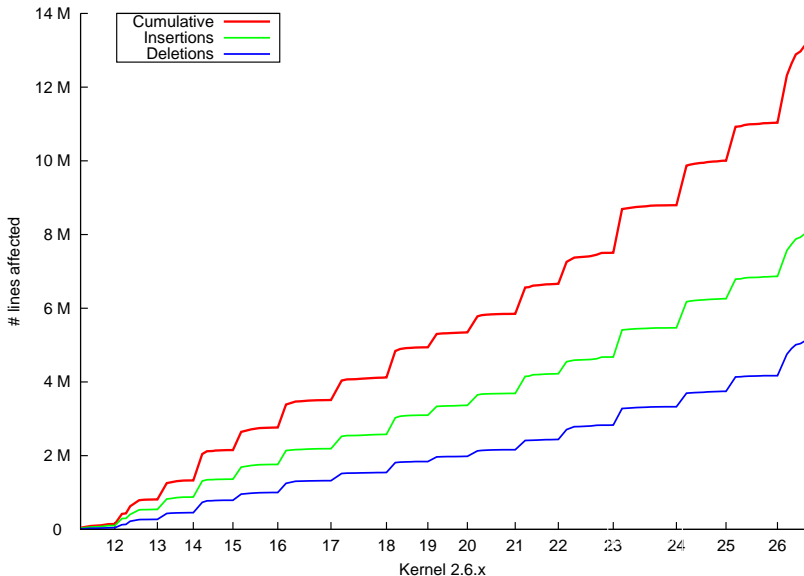
- Analysis Tools

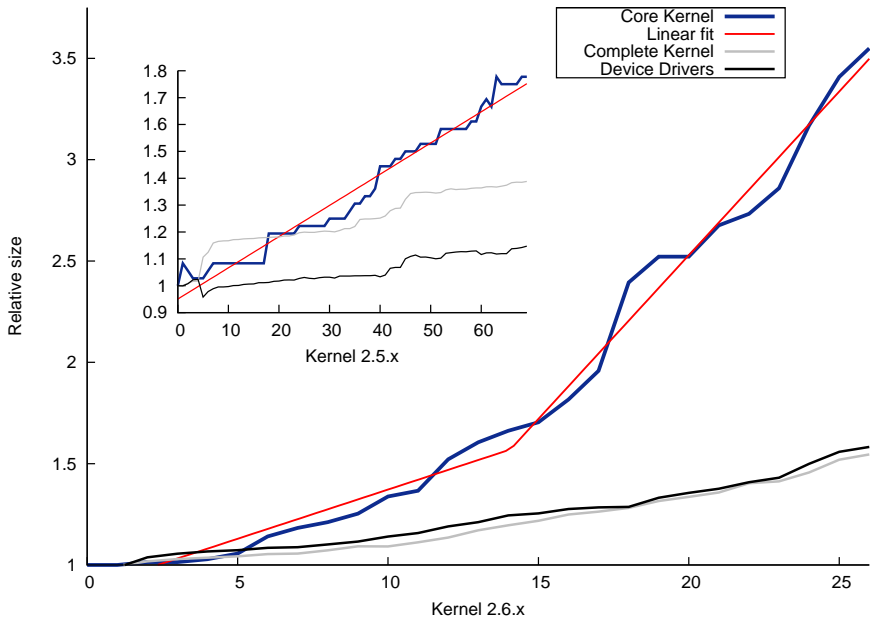
Social Aspects

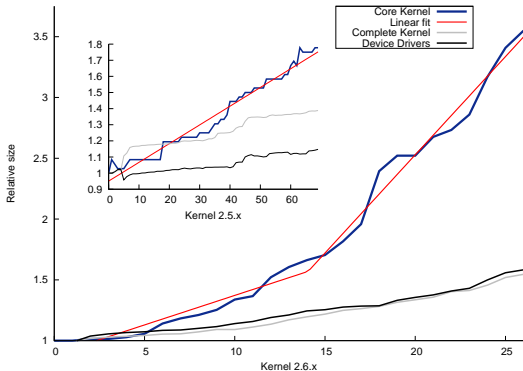
Dynamics of Kernel Development



Dynamics of Kernel Development

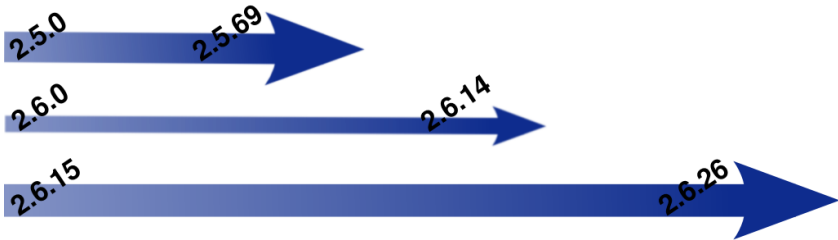


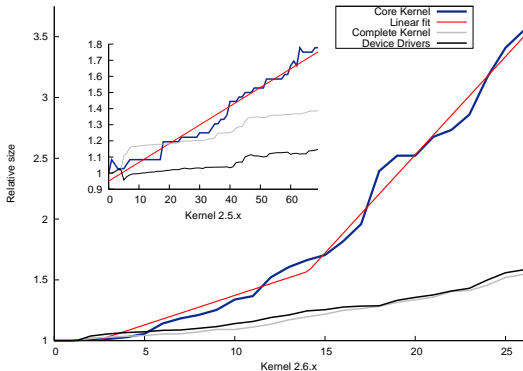




Hare and Tortoise

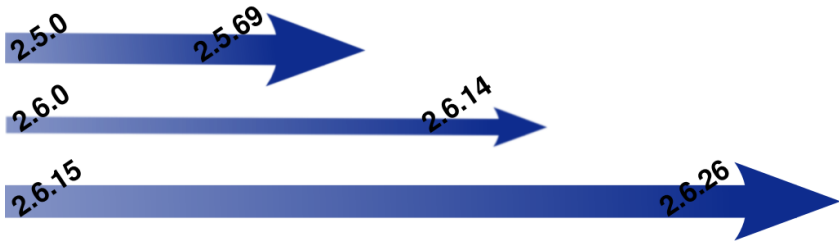
- ▶ Code sufficient?
- ▶ How to document?
- ▶ Which parts?





Hare and Tortoise

- Code sufficient?
- How to document?
- Which parts?



6 Outline

Dynamics of Kernel Development

Kernel Documentation

Understanding the Kernel

Documenting new Features

Analysis Tools

Social Aspects

7 What's available?

In-Tree

- ▶ Comments and Kernel doc
- ▶ Documentation/
- ▶ Git commit messages

External

- ▶ LKML and others
- ▶ Websites: lwn.net, kernelnewbies.org, ...
- ▶ Books and Articles

7 What's available?

In-Tree

- ▶ **Comments and Kernel doc**
- ▶ Documentation/
- ▶ Git commit messages

External

- ▶ LKML and others
- ▶ Websites: lwn.net, kernelnewbies.org, ...
- ▶ Books and Articles

```
/**
 * clocksource_khz2mult - calculates mult from khz and shift
 * @khz: Clocksource frequency in KHz
 * @shift_constant: Clocksource shift factor
 *
 * Helper functions that converts a khz counter frequency to a timsorce
 * multiplier, given the clocksource shift value
 */
static inline u32 clocksource_khz2mult(u32 khz, u32 shift_constant)
...
```

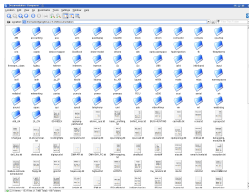
7 What's available?

In-Tree

- ▶ Comments and Kernel doc
- ▶ Documentation/
- ▶ Git commit messages

External

- ▶ LKML and others
- ▶ Websites: lwn.net, kernelnewbies.org, ...
- ▶ Books and Articles



`sx.txt - specialix SX/SI multiport serial driver readme.
Copyright (C) 1997 Roger Wolff (R.E.Wolff@BitWizard.nl)`

`...
Introduction
=====`

`This file contains some random information, that I like to have online
instead of in a manual that can get lost. Ever misplace your Linux
kernel sources? And the manual of one of the boards in your computer?
...`

7 What's available?

In-Tree

- ▶ Comments and Kernel doc
- ▶ Documentation/
- ▶ **Git commit messages**

```
commit 2087a1ad822cd3a68b73338457047fcc54da726b
Author: Gregory Haskins <ghaskins@novell.com>
Date: Fri Jun 27 14:30:00 2008 -0600
```

```
sched: add avg-overlap support to RT tasks
```

We have the notion of tracking process-coupling (a.k.a. buddy-wake) via the `p->se.last_wake` / `p->se.avg_overlap` facilities, but it is only used for cfs to cfs interactions. There is no reason why an rt to cfs interaction cannot share in establishing a relationship in a similar manner.

Because `PREEMPT_RT` runs many kernel threads as FIFO priority, we often times have heavy interaction between RT threads waking CFS applications. This patch offers a substantial boost (50-60%+) in performance under those circumstances.

External

- ▶ LKML and others
- ▶ Websites: `lwn.net`,
`kernelnewbies.org`, ...
- ▶ Books and Articles

7 What's available?

In-Tree

- ▶ Comments and Kernel doc
- ▶ Documentation/
- ▶ Git commit messages

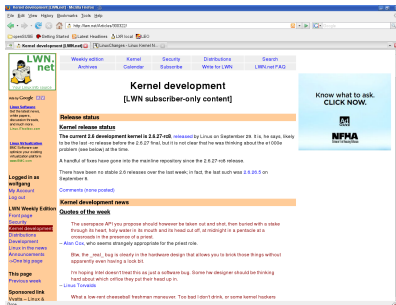
External

- ▶ **LKML and others**
- ▶ Websites: `lwn.net`,
`kernelnewbies.org`, ...
- ▶ Books and Articles

7 What's available?

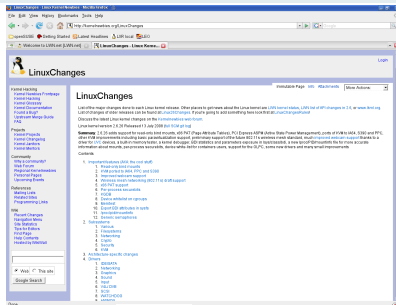
In-Tree

- Comments and Kernelsdoc
- Documentation/
- Git commit messages



External

- LKML and others
- Websites: lwn.net, kernelnewbies.org, ...
- Books and Articles



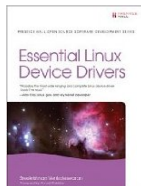
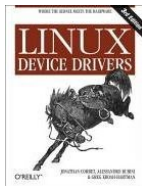
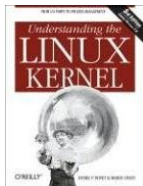
7 What's available?

In-Tree

- ▶ Comments and Kernel doc
- ▶ Documentation/
- ▶ Git commit messages

External

- ▶ LKML and others
- ▶ Websites: `lwn.net`, `kernelnewbies.org`, ...
- ▶ **Books and Articles**



7 What's available?

In-Tree

- ▶ Comments and Kernel doc
- ▶ Documentation/
- ▶ Git commit messages

External

- ▶ LKML and others
- ▶ Websites: `lwn.net`, `kernelnewbies.org`, ...
- ▶ Books and Articles

Problems

- ▶ Available? Location?
- ▶ Uptodate? Complete?

8 Summary

✓ What's good

- ▶ Huge amount of documentation available
- ▶ Implicit documentation in git
- ▶ Documentation infrastructure available

✗ What's bad

- ▶ Focus on people already intimate with the code
- ▶ Implicit documentation in git
- ▶ No consistent style
- ▶ Very fragmented and scattered

9 Outline

Dynamics of Kernel Development

Kernel Documentation

Understanding the Kernel

- Documenting new Features

- Analysis Tools

Social Aspects

Documenting new features

Completely Fair Scheduler

- ▶ Turbulent emergence
- ▶ Completely *replaces* old scheduler
- ▶ Considerable in-tree development after merge

High Resolution Timers

- ▶ Long external development
- ▶ New foundation for *existing* framework
- ▶ Merged at very mature state

Opposite strategies...

... also with respect to documentation!

Documenting new features

Completely Fair Scheduler

- ▶ Turbulent emergence
- ▶ Completely *replaces* old scheduler
- ▶ Considerable in-tree development after merge

High Resolution Timers

- ▶ Long external development
- ▶ New foundation for *existing* framework
- ▶ Merged at very mature state

Opposite strategies...

... also with respect to documentation!

Documenting new features

Completely Fair Scheduler

- ▶ Turbulent emergence
- ▶ Completely *replaces* old scheduler
- ▶ Considerable in-tree development after merge

High Resolution Timers

- ▶ Long external development
- ▶ New foundation for *existing* framework
- ▶ Merged at very mature state

Opposite strategies...

... also with respect to documentation!

Documenting new features

Completely Fair Scheduler

- ▶ Turbulent emergence
- ▶ Completely *replaces* old scheduler
- ▶ Considerable in-tree development after merge

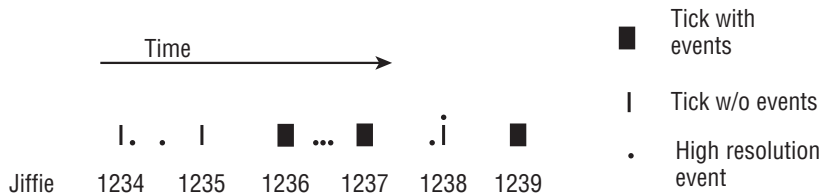
High Resolution Timers

- ▶ Long external development
- ▶ New foundation for *existing* framework
- ▶ Merged at very mature state

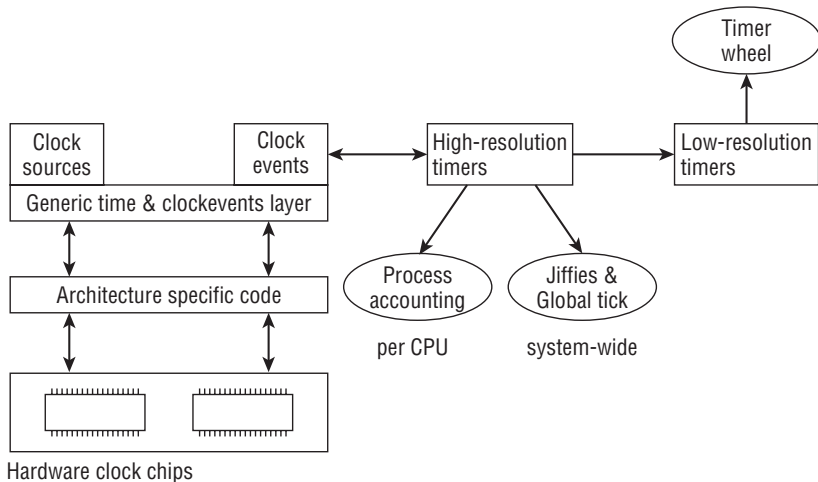
Opposite strategies. . .

. . . also with respect to documentation!

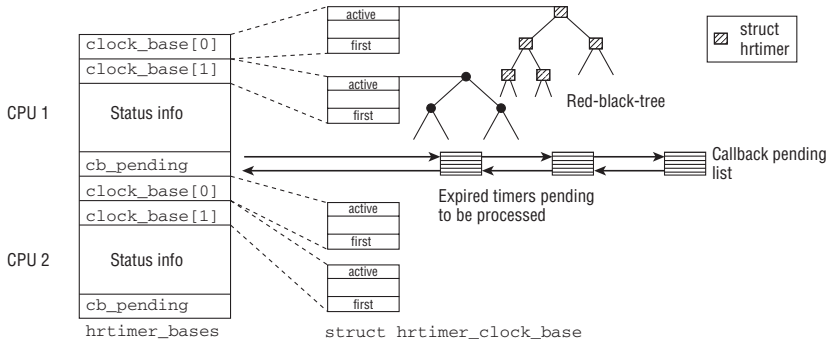
High Resolution Timers



High Resolution Timers



High Resolution Timers



```

void hrtimer_interrupt(struct clock_event_device *dev)
{
    struct hrtimer_cpu_base *cpu_base = &__get_cpu_var(hrtimer_bases);
    struct hrtimer_clock_base *base;
    ktime_t expires_next, now;
    int i, raise = 0;

    BUG_ON(!cpu_base->hres_active);
    cpu_base->nr_events++;
    dev->next_event.tv64 = KTIME_MAX;

retry:
    now = ktime_get();
    expires_next.tv64 = KTIME_MAX;
    base = cpu_base->clock_base;

    for (i = 0; i < HRTIMER_MAX_CLOCK_BASES; i++) {
        ktime_t basenow;
        struct rb_node *node;

        spin_lock(&cpu_base->lock);
        basenow = ktime_add(now, base->offset);

        while ((node = base->first)) {
            struct hrtimer *timer;

            timer = rb_entry(node, struct hrtimer, node);

            if (basenow.tv64 < timer->expires.tv64) {
                ktime_t expires;

                expires = ktime_sub(timer->expires,
                                    base->offset);
                if (expires.tv64 < expires_next.tv64)
                    expires_next = expires;
                break;
            }

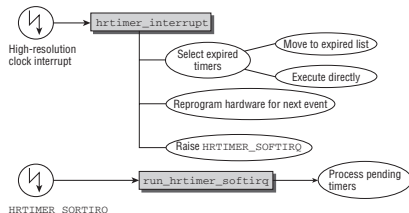
            /* Move softirq callbacks to the pending list */
            if (timer->cb_mode == HRTIMER_CB_SOFTIRQ) {
                __remove_hrtimer(timer, base,
                                HRTIMER_STATE_PENDING, 0);
                list_add_tail(&timer->cb_entry,
                              &base->cpu_base->cb_pending);
                raise = 1;
                continue;
            }

            __run_hrtimer(timer);
            spin_unlock(&cpu_base->lock);
            base++;
        }
        cpu_base->expires_next = expires_next;

        /* Reprogramming necessary ? */
        if (expires_next.tv64 != KTIME_MAX) {
            if (!tick_program_event(expires_next, 0))
                goto retry;
        }

        /* Raise softirq ? */
        if (raise)
            raise_softirq(HRTIMER_SOFTIRQ);
    }
}

```



High Resolution Timers

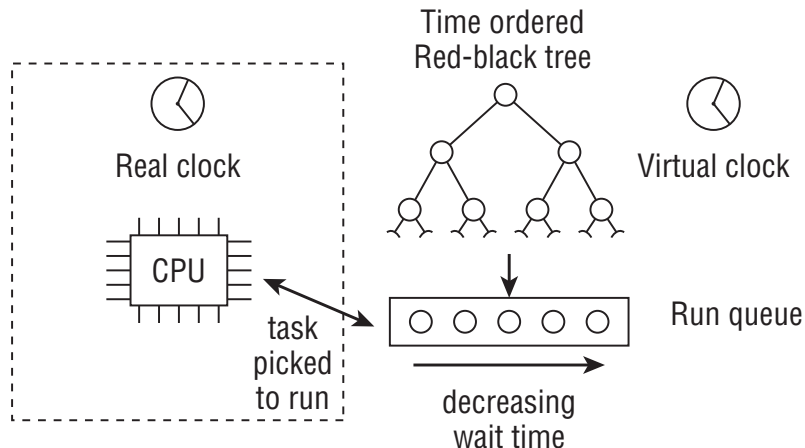
Available

- ▶ Orthogonal patch structure
- ▶ Component submission
- ▶ Design Documentation

Challenges

- ▶ Introduce conceptual parts
- ▶ Disentangle alternatives
- ▶ Prioritise important against unimportant code

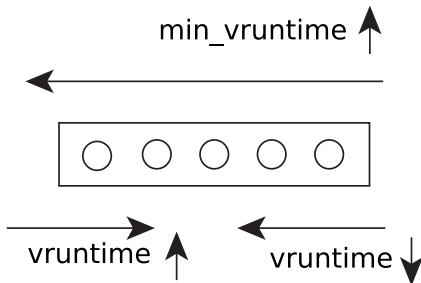
Completely Fair Scheduler



Completely Fair Scheduler

Virtual clock replacement

```
static inline s64 entity_key(struct cfs_rq *cfs_rq,  
                             struct sched_entity *se) {  
    return se->vruntime - cfs_rq->min_vruntime;  
}
```



Conclusions

- ▶ Code history can ease documentation
- ▶ Identify stable components
- ▶ Reduce complexity

16 Analysis Tools



16 Analysis Tools



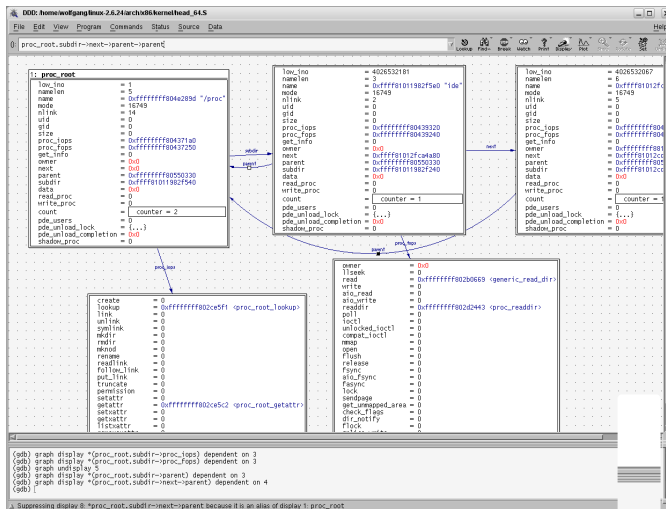
LXR Source Code Cross Reference

```

4141 __linkage void __sched_schedule(void)
4142 {
4143     struct task_struct *prev, *next;
4144     unsigned long *switch_count;
4145     struct rq *rq;
4146     int cpu;
4147
4148     need_resched:
4149     preempt_disable();
4150     cpu = smp_processor_id();
4151     rq = cpu_rq(cpu);
4152     rcu_qsctr_inc(cpu);
4153     prev = rq->curr;
4154     switch_count = &prev->nivcsw;
4155
4156     release_kernel_lock(prev);
4157     need_resched_nonpreemptible:
4158
4159     schedule_debug(prev);
4160     hrtick_clear(rq);
4161
4162     /*
4163     * Do the rq-clock update outside the rq lock:
4164     */
4165     local_irq_disable();
4166     update_rq_clock(rq);
4167     spin_lock(&rq->lock);
4168     clear_tsk_need_resched(prev);
4169
4170     if (prev->state && !(preempt_count() & PREEMPT_ACTIVE)) {
4171         if (unlikely(signal_pending_state(prev->state, prev)))
4172             prev->state = TASK_RUNNING;
4173         else
4174             deactivate_task(rq, prev, 1);
4175         switch_count = &prev->nivcsw;
4176     }
4177
4178 #ifdef CONFIG_SMP
4179     if (prev->sched_class->pre_schedule)
4180         prev->sched_class->pre_schedule(rq, prev);
4181 #endif
4182
4183     if (unlikely(!rq->nr_running))
4184         idle_balance(cpu, rq);
4185
4186     prev->sched_class->put_prev_task(rq, prev);
4187
4188     Done

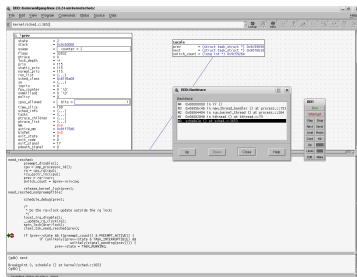
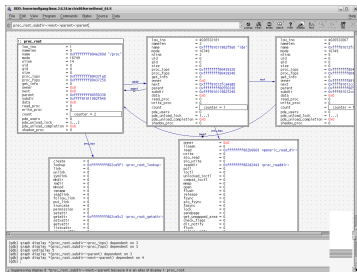
```


18 (K)GDB and DDD



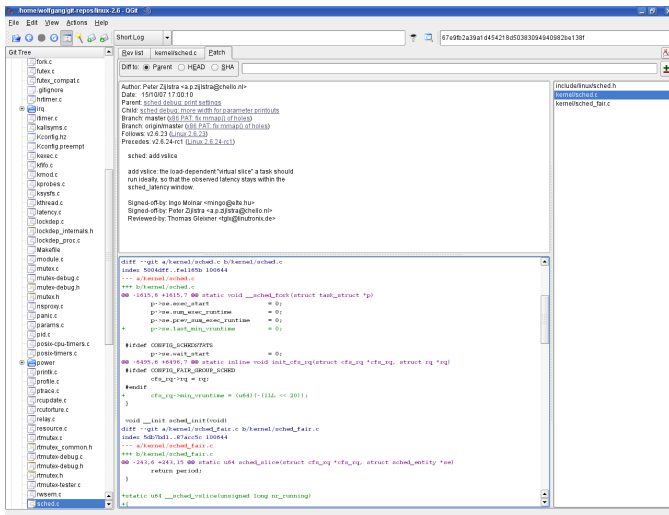


18 (K)GDB and DDD

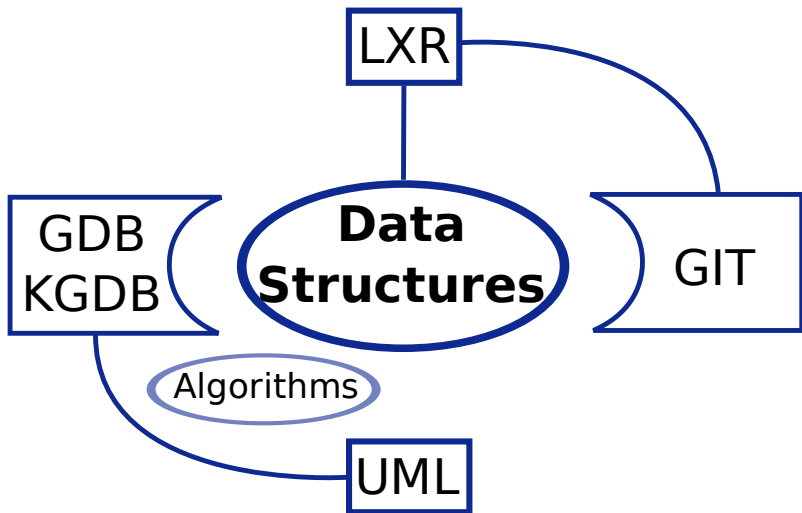


Andrew Morton on KGDB

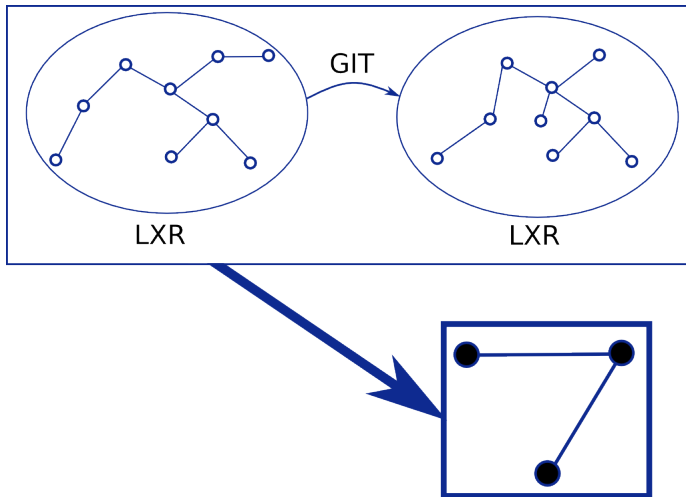
I used kgdb continuously for 4-5 years until it broke. I don't think I ever used it much for “debugging” as such. I used it more for general observation of what’s going on in the kernel. And for *confirmation* of what’s going on (ie: testing that the actual state matches the expected state).



How it all fits together



How it all fits together



Dynamics of Kernel Development

Kernel Documentation

Understanding the Kernel

Documenting new Features

Analysis Tools

Social Aspects

From: Con Kolivas

Ooh you have a vm patch that helps swap on the desktop! I can help you here with my experience from swap prefetch.

1. Get it reviewed and have noone show any evidence it harms
2. Find hundreds of users who can testify it helps
3. Find a way of quantifying it.
4. ...
5. Merge into mainline.

There, that should get you as far as 4. I haven't figured out what 4 is yet. I believe it may be goto 1;

From: Con Kolivas

Ooh you have a vm patch that helps swap on the desktop! I can help you here with my experience from swap prefetch.

1. Get it reviewed and have noone show any evidence it harms
2. Find hundreds of users who can testify it helps
3. Find a way of quantifying it.
4. ...
5. Merge into mainline.

There, that should get you as far as 4. I haven't figured out what 4 is yet. I believe it may be goto 1;

Linus on smelly pets

Ok, so now that I've insulted you and your pets (they're ugly!), show me wrong, and then call me a d*ckhead.
("Linus - you're a d*ckhead, and you didn't understand the problem, so you're a *stupid* d*ckhead. And my pet may be ugly, but yours *smells* bad!").

From: Rusty Russel

```
-#define ARRAY_SIZE(x) (sizeof(x) / sizeof((x)[0]))  
+#define ARRAY_SIZE(arr) (sizeof(arr) / sizeof((arr)[0]) \br/>+ + sizeof(typeof(int[1 - 2*!!__builtin_types_compatible_p(typeof(arr), \br/>+   typeof(&arr[0]))]))*0)
```

Reply from Linus Torvalds

Rusty, that's a work of art.

However, I would suggest that you never show it to anybody ever again. I'm sure that in fifty years, it will be worth much more. So please keep it tightly under wraps, to keep people from gouging their eyes out~W~W~W~W~W~W make a killing in the art market.

Improved patch

OK, many people complained that it needed a comment. Good point!

==

Add comment to ARRAY_SIZE macro.

diff -r 933e410f204f include/linux/kernel.h

--- a/include/linux/kernel.h Sat Mar 10 09:55:31 2007 +1100

+++ b/include/linux/kernel.h Sat Mar 10 09:55:53 2007 +1100

<at> <at> -35,6 +35,7 <at> <at> extern const char linux_proc_banner[];

#define ALIGN(x,a) __ALIGN_MASK(x,(typeof(x))(a)-1)

#define __ALIGN_MASK(x,mask) (((x)+(mask))&~(mask))

/* GCC is awesome. */

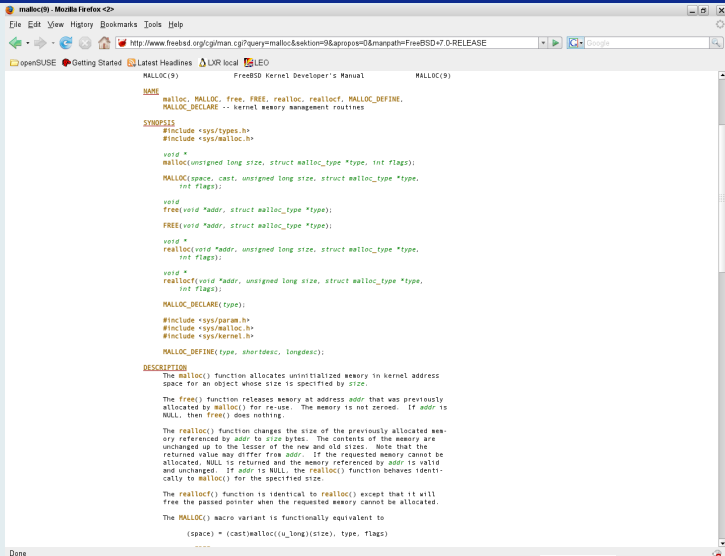
#define ARRAY_SIZE(arr) (sizeof(arr) / sizeof((arr)[0]) \
+ sizeof(typeof(int[1 - 2*!!__builtin_types_compatible_p(typeof(arr), \
typeof(&arr[0]))]))*0)

Thanks for you attention!

A Comparison: Allocating Memory

Everybody needs memory

- ▶ Core OS service
- ▶ Stable interface (introduced \approx v0.98)
- ▶ Documentation situation representative



www.freebsd.org/cgi/man.cgi?query=malloc&apropos=0&sektion=0&manpath=FreeBSD+7.0-RELEASE&format=html

I/O Kit Device Driver Design Guidelines



Introduction

The Ilibern C++ Runtime

- Creation of the Runtime System
- Object Creation and Destruction
- Object Inspection and Dynamic Casting
- Binary Compatibility

Ilibern Collection and Container Classes

- The Ilibern Classes and Your Driver
- Ilibern Collection and Container Class Overview
- Using the Ilibern Collection and Container Classes
- Configuring Your Driver Using the Ilibern Classes

The IOService API

Making Hardware Accessible to Applications

Kernel-User Notification

Displaying Localized Information About Drivers

- Internationalizing Kernel Extensions
- Getting the Path to a KEXT From User Space

Debugging Drivers

Testing and Deploying Drivers

- Testing Strategies
- Packaging Drivers for Installation

Developing a Device Driver to Run on an Intel-Based Macintosh

- Glossary
- Revision History
- Index



Developer Connection



Log In | Not a Member?

ADC Home > Reference Library > Guides > Hardware & Drivers > I/O Kit Device Driver Design Guidelines > The Ilibern C++ Runtime >

Hide TOC

Advanced Search



Contact ADC

< Previous Page Next Page >

Object Creation and Destruction

Because exceptions are excluded from the kernel's restricted form of C++, you cannot implement "normal" C++ constructors and destructors without jeopardy. Constructors and destructors are typed to return no value (such as an error code). Normally, if they encounter a problem, they raise an exception. But because exceptions aren't supported in the kernel's C++ runtime, there is no way for you to know when an allocation or deallocation error has occurred.

This situation prompted a design feature of the Ilibern's C++ runtime system that uses OSMetaClass macros to specify the structure of a class—that is, the metaclass data structures and functional interfaces—for the runtime typing system. The macros also define the primary constructor and a destructor for a class. These macro-created constructors are guaranteed not to fail because they do not themselves perform any allocations. Instead, the runtime system defers the actual allocation of objects until their initialization (usually in the `init` member function). Because the `init` function is typed to return a `bool`, it makes it possible to return an error upon any failure.

In this section:

Using the OSMetaClass Constructor Macros
Allocating Objects Dynamically
Global Initializers

Using the OSMetaClass Constructor Macros

When you create a C++ class based on OSObject, your code must call a matching pair of macros based upon the OSMetaClass class. The calls must be among the first statements in both the definition and implementation of the class. These macros are critical to your class because they enter metaclass information about it into the Ilibern runtime typing facility and define the static constructor and destructor for your class.

For concrete (that is, non-abstract) classes, the first macro, `OSDeclareDefaultConstructors`, declares the C++ constructors; by convention you insert this macro as the first element of the class declaration in the header file. For example:

```
class com_myCompany_driver_MyDriver : public IOService
{
    OSDeclareDefaultConstructors(com_myCompany_driver_MyDriver);
    /* ... */
};
```

Your class implementation must include the companion "define" macro, `OSDefineMetaClassAndConstructors`. This macro defines the constructor and destructor, implements the OSMetaClass allocation member function (`alloc`) for the class, and supplies the metaclass information for the runtime typing system. `OSDefineMetaClassAndConstructors` takes as arguments the name of your driver's class and the name of its base class. It creates these in namespace code that allows your driver class to be loaded and instantiated while the kernel is

developer.apple.com/documentation/DeviceDrivers/Conceptual/WritingDeviceDriver/CPlusPlusRuntime/chapter_2_section_3.html#/apple_ref/doc/uid/TP30000695-BAJCCBGJ

ExAllocatePoolWithTag - Mozilla Firefox <2>

File Edit View History Bookmarks Tools Help

http://msdn.microsoft.com/en-us/library/ms796989.aspx

openSUSE Getting Started Latest Headlines UXR local LEO

United States - English Microsoft.com Welcome Sign in

msdn Search MSDN with Live Search Web MSDN Home Developer Centers

Microsoft Developer Network

Home Library Learn Downloads Support Community

Printer Friendly Version Send Add Content... Click to Rate and Give Feedback

Getting Started with Windows Drivers

Kernel-Mode Driver Architecture

Design Guide

Reference

Standard Driver Routines

Driver Support Routines

Summary of Kernel-Mode

Object Manager Routines

Memory Manager Routines

Process and Thread Manag

I/O Manager Routines

Power Manager Routines

Configuration Manager Ro

Kernel Transaction Manag

Security Reference Monitor

Core Kernel Library Suppo

Executive Library Support

ExAcquireFastMutex

ExAcquireFastMutexUn

ExAcquireResourceExc

ExAcquireResourceExc

ExAcquireResourceSh

ExAcquireResourceSh

ExAcquireSharedStar

ExAcquireSharedWaitF

ExAllocateFromNPaged

ExAllocateFromPagedL

ExAllocateFromZone

ExAllocatePool

ExAllocatePoolWithQuo

ExAllocatePoolWithQuo

ExAllocatePoolWithTag

ExAllocatePoolWithTag

Windows Driver Kit: Kernel-Mode Driver Architecture

ExAllocatePoolWithTag

The **ExAllocatePoolWithTag** routine allocates pool memory of the specified type and returns a pointer to the allocated block.

```

NTSTATUS
ExAllocatePoolWithTag(
    IN POOL_TYPE PoolType,
    IN SIZE_T NumberOfBytes,
    IN ULONG Tag
);
    
```

Parameters

PoolType
Specifies the type of pool memory to allocate. Each allocation code path should use a unique pool tag to help debuggers and verifiers identify the code path. You can modify the PoolType value by using a bitwise OR with the POOL_RAISE_IF_ALLOCATION_FAILURE flag. This flag causes an exception to be raised if the request cannot be satisfied. Similarly, you can modify PoolType by using a bitwise OR with the POOL_COLD_ALLOCATION flag as a hint to the kernel to allocate the memory from pages that are likely to be paged out quickly. To reduce the amount of resident pool memory as much as possible, you should not reference these allocations frequently. The POOL_COLD_ALLOCATION flag is only advisory and is available for Microsoft Windows XP and later operating systems. For a description of the available pool memory types, see [POOL_TYPE](#).

NumberOfBytes
Specifies the number of bytes to allocate.

Tag
Specifies the pool tag for the allocated memory. Specify the pool tag as a character literal of up to four characters delimited by single quotation marks (for example, Tag1). The string is usually specified in reverse order (for example, "1gaT"). The ASCII value of each character in the tag must be between 0 and 127. Every allocation code path should use a unique pool tag to ensure that debuggers and verifiers identify a distinct allocated block.

Return Value

ExAllocatePoolWithTag returns NULL if there is insufficient memory in the free pool to satisfy the request. Use of POOL_RAISE_IF_ALLOCATION_FAILURE is not recommended because it is costly. Any successful allocation that requests `NumberOfBytes > PAGE_SIZE` wastes all unused bytes on the last-allocated page. However, on Windows Vista and later, the unused bytes are no longer wasted.

Comments

This routine is used for the general pool allocation of memory.

If `NumberOfBytes` is `PAGE_SIZE` or greater, a page-aligned buffer is allocated. Memory allocations of `PAGE_SIZE` or less do not cross page boundaries. Memory allocations of less than `PAGE_SIZE` are not necessarily page-aligned but are aligned on an 8-byte boundary.

msdn.microsoft.com/en-us/library/ms796989.aspx

PERLGUTS (1) Perl Programmers Reference Guide PERLGUTS (1)

Memory Allocation

It is suggested that you use the version of malloc that is distributed with Perl. It keeps pools of various sizes of unallocated memory in order to satisfy allocation requests more quickly. However, on some platforms, it may cause spurious malloc or free errors.

```
New(x, pointer, number, type);
Newc(x, pointer, number, type, cast);
Newx(x, pointer, number, type);
```

These three macros are used to initially allocate memory.

The first argument *x* was a "magic cookie" that was used to keep track of who called the macro, to help when debugging memory problems. However, the current code makes no use of this feature (most Perl developers now use run-time memory checks), so this argument can be any number.

The second argument *pointer* should be the name of a variable that will point to the newly allocated memory.

The third and fourth arguments *number* and *type* specify how many of the specified type of data structure should be allocated. The argument *type* is passed to `sizeof`. The final argument to `Newc`, *cast*, should be used if the *pointer* argument is different from the *type* argument.

Unlike the `New` and `Newc` macros, the `Newx` macro calls `memzero` to zero out all the newly allocated memory.

```
Renew(pointer, number, type);
Renewc(pointer, number, type, cast);
SafeFree(pointer);
```

These three macros are used to change a memory buffer size or to free a piece of memory no longer needed. The arguments to `Renew` and `Renewc` match those of `New` and `Newc` with the exception of not needing the "magic cookie" argument.

```
Move(source, dest, number, type);
Copy(source, dest, number, type);
Zero(dest, number, type);
```

These three macros are used to move, copy, or zero out previously allocated memory. The *source* and *dest* arguments point to the source and destination starting points. Perl will move, copy, or zero out *number* instances of the size of the *type* data structure (using the `sizeof` function).

PerlIO

The most recent development releases of Perl have been experimenting with removing Perl's dependency on the "normal" standard I/O suite and allowing other stdio implementations to be used. This involves creating a new abstraction layer that then calls whichever implementation of stdio Perl was compiled with. All XS/IOs should now use the functions in the PerlIO abstraction layer and not make any assumptions about what kind of stdio is being used.

2nd Berkeley Distribution perl 5.005, patch 02 PERLGUTS (1) - 21

Chapter 5. Memory Management in Linux

Table of Contents

[The Slab Cache](#)

[User Space Memory Access](#)

[More Memory Management Functions](#)

The Slab Cache

[kcalloc](#) — allocate memory for an array. The memory is set to zero.

[kmallocc_node](#) — allocate memory from a specific node

[kzalloc](#) — allocate memory. The memory is set to zero.

[kzalloc_node](#) — allocate zeroed memory from a particular memory node.

[kmemp_cache_create](#) — Create a cache.

[kmemp_cache_shrink](#) — Shrink a cache.

[kmemp_cache_destroy](#) — delete a cache

[kmemp_cache_alloc](#) — Allocate an object

[kmemp_cache_free](#) — Deallocate an object

[kfree](#) — free previously allocated memory

Name

`kmalloc_node` — allocate memory from a specific node

Synopsis

```
void * kmalloc_node (size_t size,  
                    gfp_t flags,  
                    int node);
```

Arguments

size

how many bytes of memory are required.

flags

the type of memory to allocate (see `kcalloc`).

node

node to allocate from.

Description

`kcalloc` for non-local nodes, used to allocate from a specific node if available. Equivalent to `kalloc` in the non-NUMA single-node case.

Name

kcalloc — allocate memory for an array. The memory is set to zero.

Synopsis

```
void * kcalloc (size_t n,  
               size_t size,  
               gfp_t flags);
```

Arguments

n
number of elements.

size
element size.

flags
the type of memory to allocate.

Description

The *flags* argument may be one of:

GFP_USER - Allocate memory on behalf of user. May sleep.

GFP_KERNEL - Allocate normal kernel ram. May sleep.

GFP_ATOMIC - Allocation will not sleep. May use emergency pools. For example, use this inside interrupt handlers.

GFP_HIGHUSER - Allocate pages from high memory.

GFP_NOIO - Do not do any I/O at all while tries to get memory.

www.kernel.org/doc/htmldocs/kernel-api/ch05.html

That's what you also get. . .

