

# **Séminaire « Linux et le temps réel » 5 mai 2010**

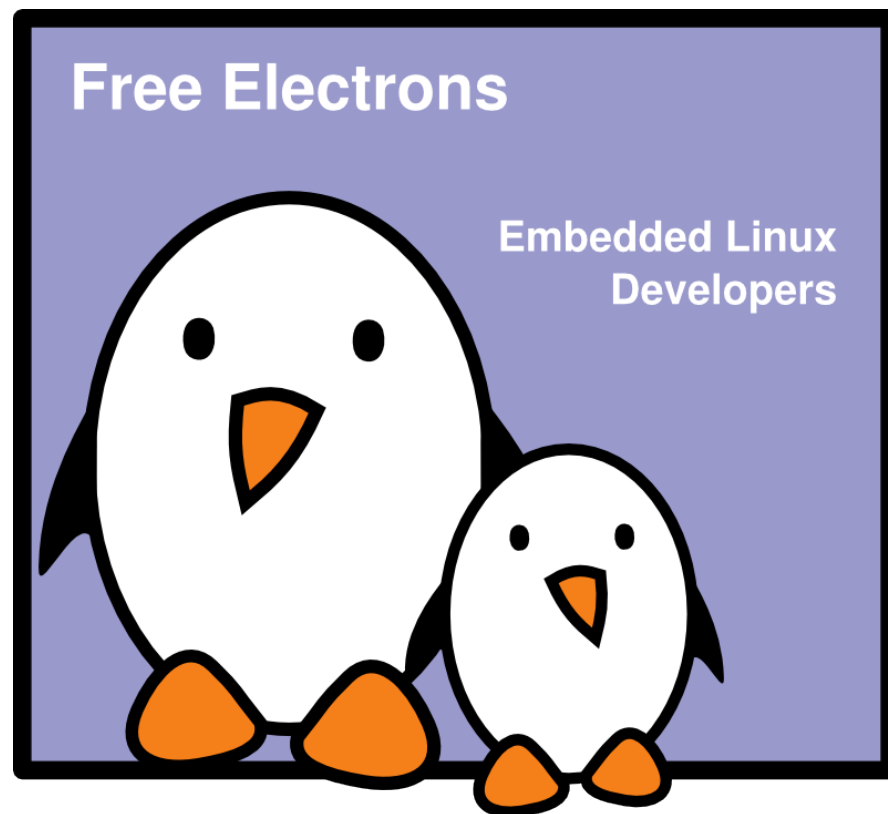
*animé par*  
Thomas Petazzoni, Bootlin  
<https://bootlin.com>

*organisé par*  
Captronic en collaboration avec le Pôle Aerospace  
Valley et Midi Pyrénées Innovation



## Linux and real-time

Michael Opdenacker  
Thomas Petazzoni  
**Bootlin**



© Copyright 2004-2010, Bootlin.  
Creative Commons BY-SA 3.0 license  
Latest update: Jul 27, 2018,  
Document sources, updates and translations:  
<https://bootlin.com/doc/training/embedded-linux/>  
Corrections, suggestions, contributions and translations are welcome!



# Bootlin (1)

- ▶ Embedded Linux experts

- ▶ Expertise in the low-level Linux stack : bootloader, kernel, libraries, applications
- ▶ Founded in 2004
- ▶ Strong link with the open-source community
- ▶ Located in Nice and Toulouse

- ▶ Development services

- ▶ Kernel development and drivers : BSP development, porting existing BSP to recent kernel versions, mainstreaming, etc.
- ▶ Embedded Linux system integration : building root filesystems, integration of useful open-source components for your system
- ▶ System optimization : boot time, power management, etc.
- ▶ Open source components customization



# Bootlin (2)

- ▶ Partner with Calao Systems for hardware development
  - ▶ AT91 and OMAP expertise
- ▶ Training services
  - ▶ *Embedded Linux system development* course (5 days)
  - ▶ *Linux kernel driver development* course (5 days)
  - ▶ On-site and public training sessions
  - ▶ All materials freely published online
- ▶ Customers: Texas Instruments, Windriver, Motorola, Freescale, Alstom, Nokia Siemens Networks, Ikusi, CS, Chess, Logiplus, etc.
- ▶ <https://bootlin.com/>



# Thomas Petazzoni

- ▶ At Bootlin since 2008
  - ▶ Embedded Linux development: power management projects, boot time optimizations projects, system integration projects, etc.
  - ▶ More than 100 days of training around the world.
- ▶ Open-source contributions
  - ▶ Major contributor to Buildroot, an embedded Linux build system
  - ▶ Contributions to the Linux kernel, Qemu, U-Boot, etc.
- ▶ In the past, 3.5 years experience with the Linux kernel: storage virtualization driver and port to a MIPS platform
- ▶ 10+ years Linux user and developer.



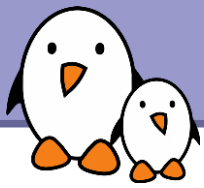
# Agenda : morning

- ▶ Embedded Linux : introduction
- ▶ Linux and real-time : history, issues and solutions
- ▶ Focus on PREEMPT\_RT
  - ▶ How it works
  - ▶ How to set it up
  - ▶ How to develop applications for PREEMPT\_RT
- ▶ Focus on Xenomai
  - ▶ How it works
  - ▶ How to set it up
  - ▶ How to develop applications for PREEMPT\_RT



# Agenda : afternoon

- ▶ Pratical experimentations with boards based on the ARM Atmel AT91SAM9263 processor, clocked at 180 Mhz, 64 MB RAM
  - ▶ Test scheduling latencies with a normal Linux kernel
  - ▶ Set up the PREEMPT\_RT solution
  - ▶ Test scheduling latencies with PREEMPT\_RT
  - ▶ Set up the Xenomai solution
  - ▶ Test scheduling latencies with Xenomai



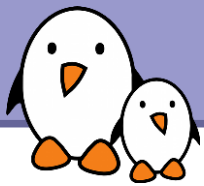
# Linux and real-time





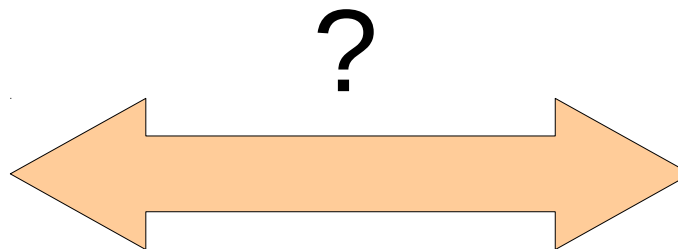
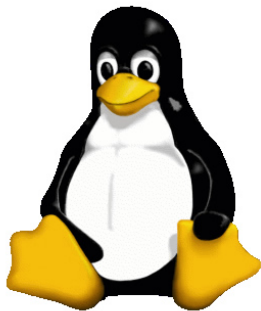
# Embedded Linux

- ▶ The Free Software and Open Source world offers a broad range of tools to develop embedded systems.
- ▶ Advantages
  - ▶ Reuse of existing components for the base system.  
Allows to focus on the added value of the product.
  - ▶ High quality, proven components (Linux kernel, C libraries...)
  - ▶ Complete control on the choice of components.  
Modifications possible without external constraints.
  - ▶ Community support: tutorials, mailing lists...
  - ▶ Low cost, in particular no per-unit royalties.
  - ▶ Potentially less legal issues.
  - ▶ Easier access to software and tools.



# Embedded Linux and real time

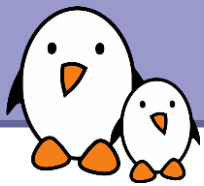
- ▶ Due to its advantages, Linux and the open-source softwares are more and more commonly used in embedded applications
- ▶ However, some applications also have real-time constraints
- ▶ They, at the same time, want to
  - ▶ Get all the nice advantages of Linux : hardware support, components re-use, low cost, etc.
  - ▶ Get their real-time constraints met



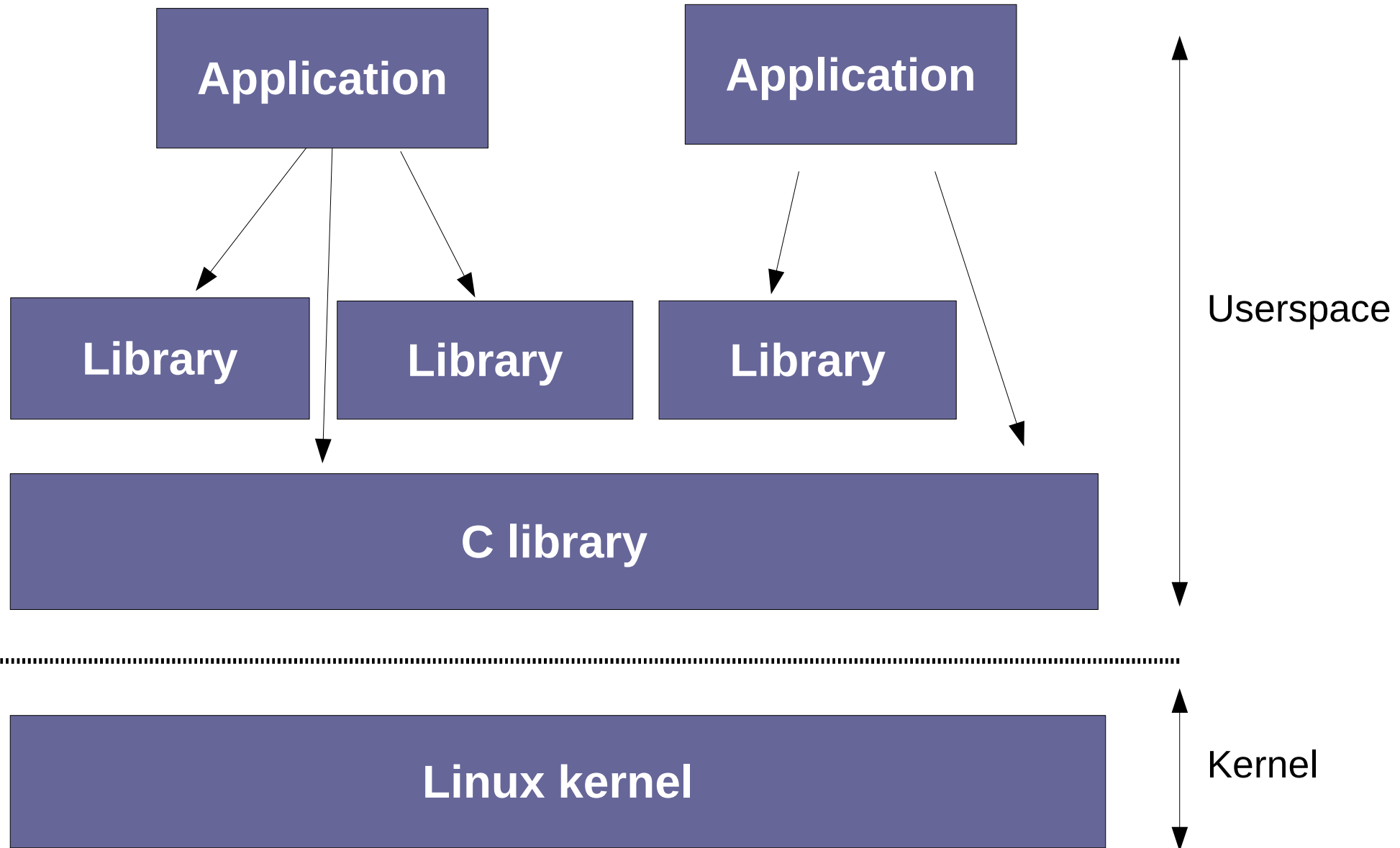


# Embedded Linux and real time

- ▶ Linux is an operating system part of the large Unix family
- ▶ It was originally designed as a time-sharing system
  - ▶ The main goal is to get the best throughput from the available hardware, by making the best possible usage of resources (CPU, memory, I/O)
  - ▶ Time determinism is not taken into account
- ▶ On the opposite, real-time constraints imply time determinism, even at the expense of lower global throughput
- ▶ Best throughput and time determinism are contradictory requirements



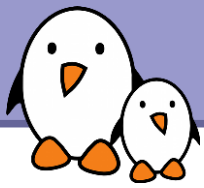
# Global architecture of a Linux system





# Linux and real-time approaches

- ▶ Over time, two major approaches have been taken to bring real-time requirements into Linux
- ▶ **Approach 1**
  - ▶ Improve the Linux kernel itself so that it matches real-time requirements, by providing bounded latencies, real-time APIs, etc.
  - ▶ Approach taken by the PREEMPT\_RT project.
- ▶ **Approach 2**
  - ▶ Add a layer below the Linux kernel that will handle all the real-time requirements, so that the behaviour of Linux doesn't affect real-time tasks.
  - ▶ Approach taken by Xenomai, RTAI and RTLinux.



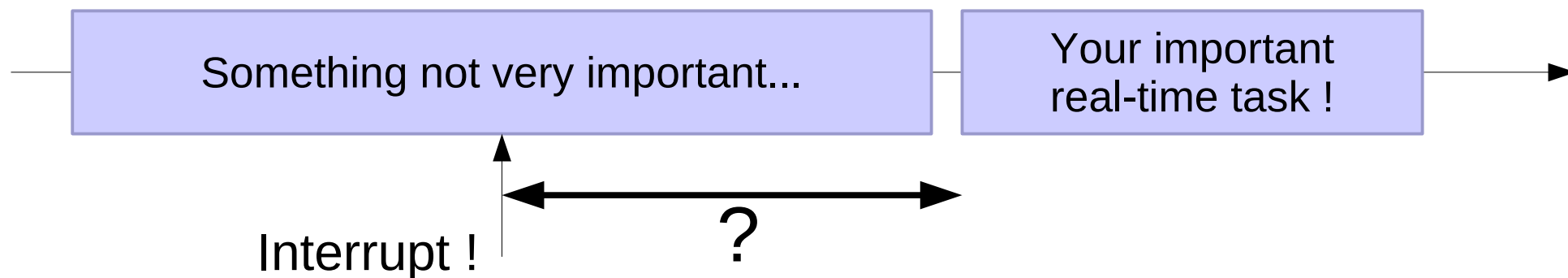
# Approach 1

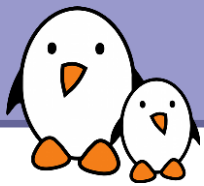
## PREEMPT\_RT



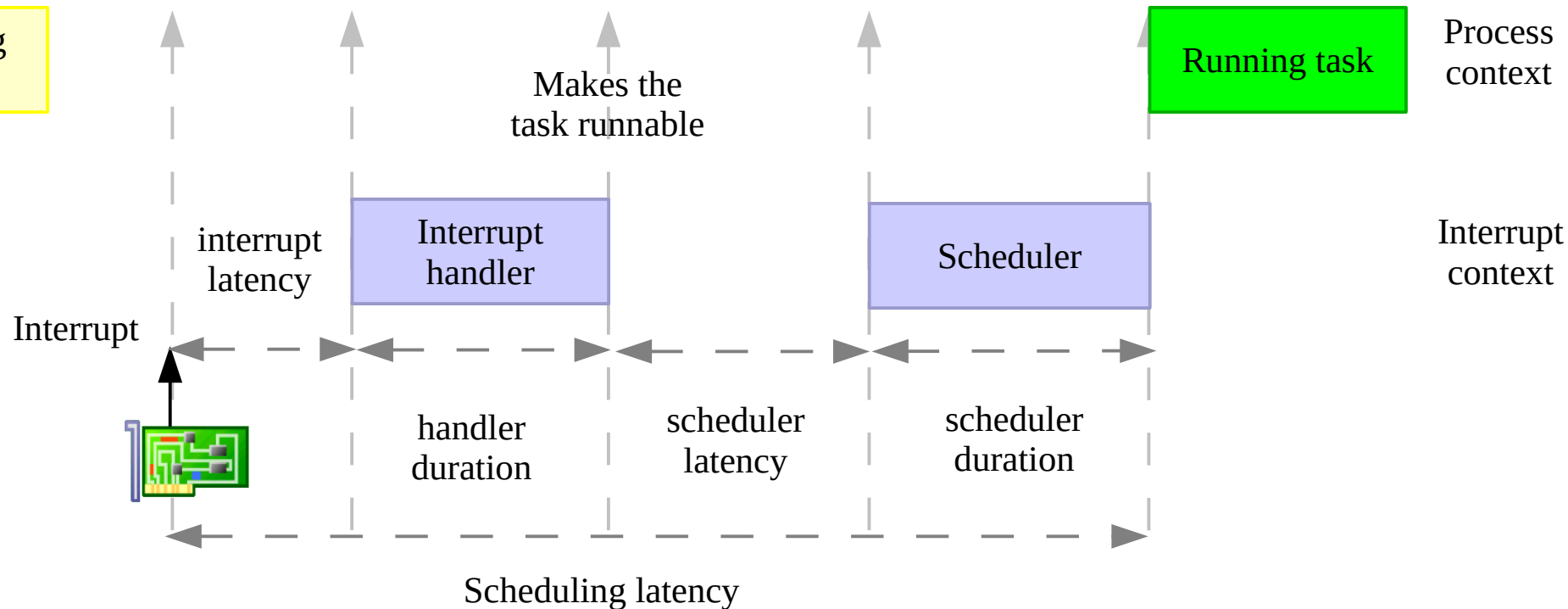
# Understanding latency

- ▶ When developing real-time applications with a system such as Linux, the typical scenario is the following
  - ▶ An event from the physical world happens and gets notified to the CPU by means of an interrupt
  - ▶ The interrupt handler recognizes and handles the event, and then wake-up the user-space task that will react to this event
  - ▶ Some time later, the user-space task will run and be able to react to the physical world event
- ▶ Real-time is about providing guaranteed worst case latencies for this reaction time, called *latency*



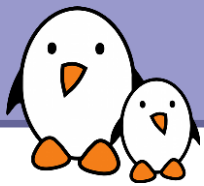


# Linux kernel latency components

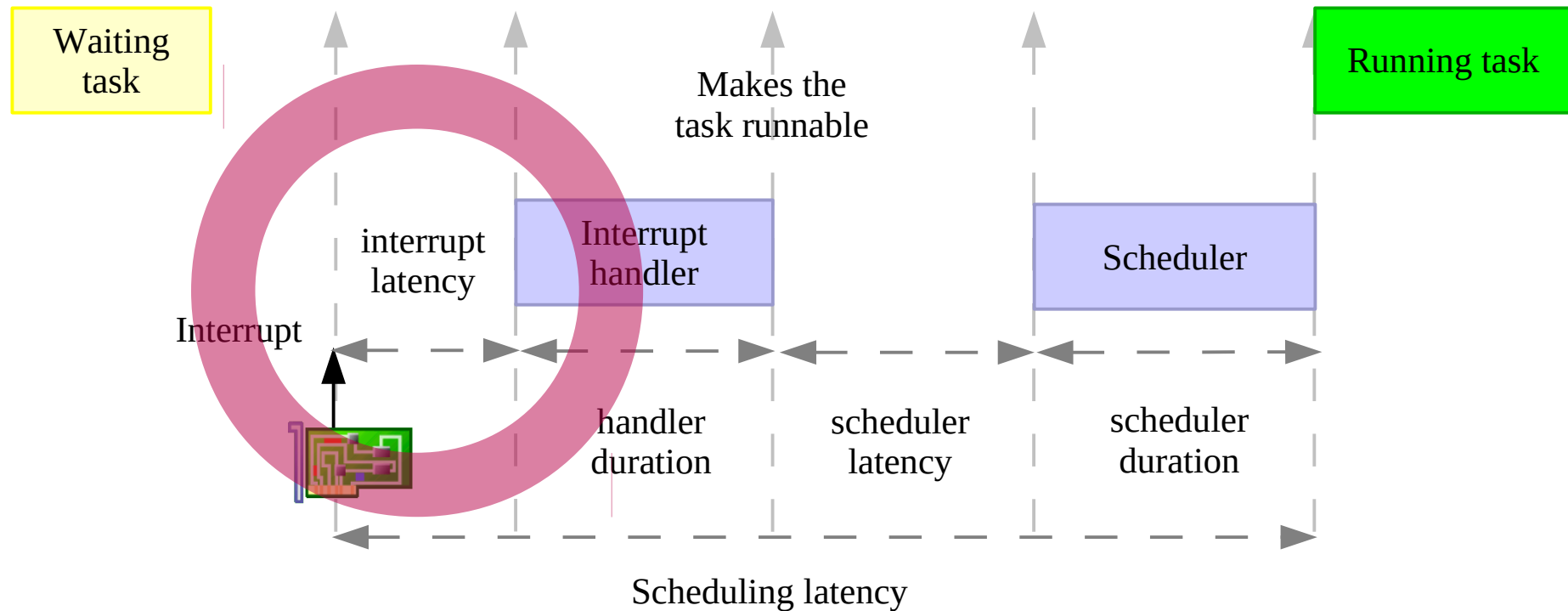


kernel latency = interrupt latency + handler duration  
+ scheduler latency + scheduler duration





# Interrupt latency

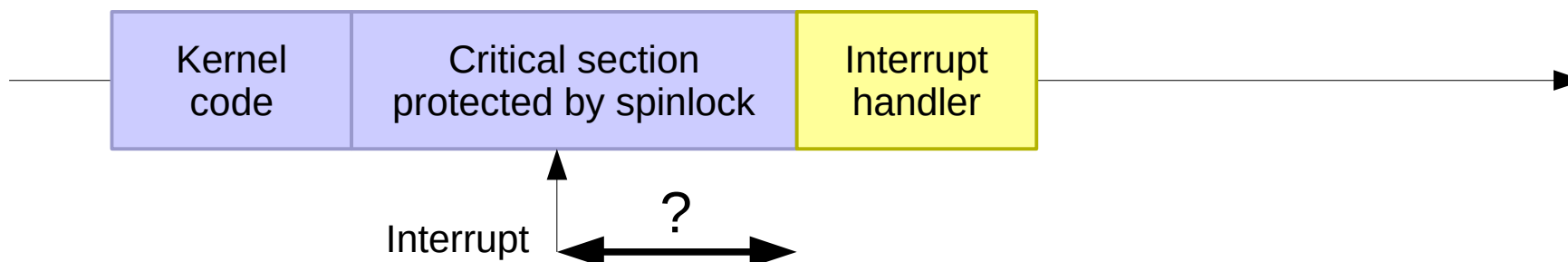


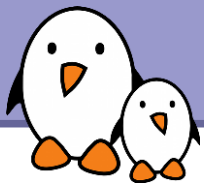
Time elapsed before executing the interrupt handler



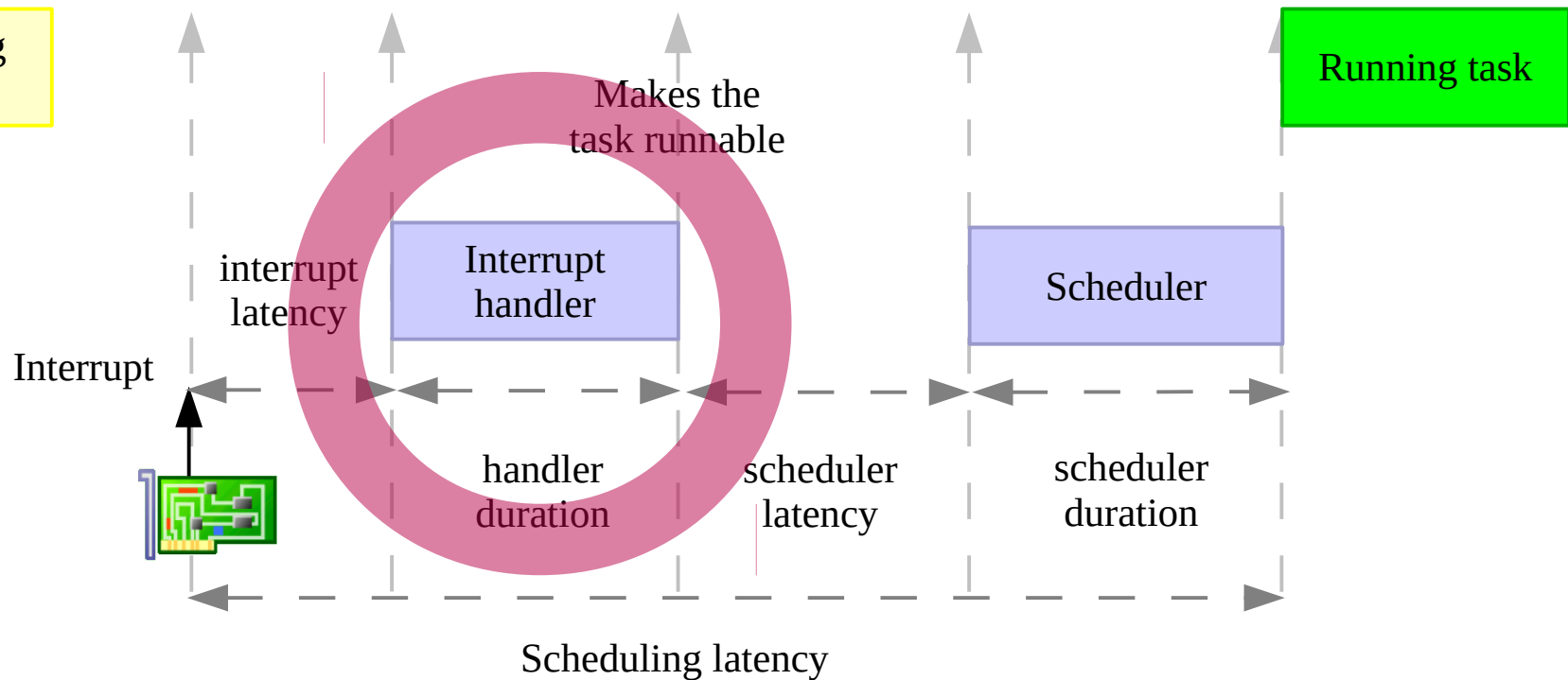
# Source of interrupt latency

- ▶ One of the concurrency prevention mechanism used in the kernel is the **spinlock**
- ▶ It has several variants, but one of the variant commonly used to prevent concurrent accesses between a process context and an interrupt context works by disabling interrupts
- ▶ Critical sections protected by spinlocks, or other section in which interrupts are explicitly disabled will delay the beginning of the execution of the interrupt handler
  - ▶ The duration of these critical sections is unbounded
- ▶ Other possible source: shared interrupts

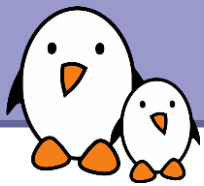




# Interrupt handler duration

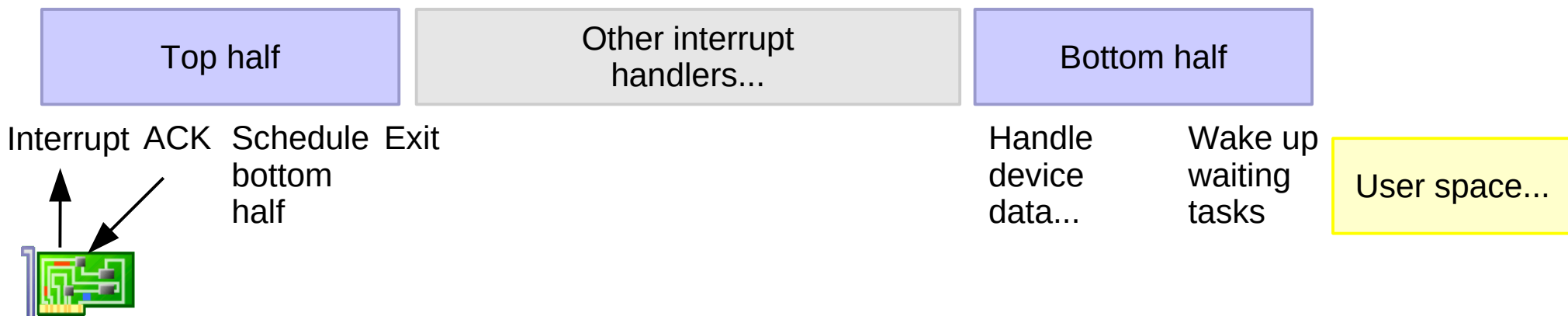


Time taken to execute the interrupt handler



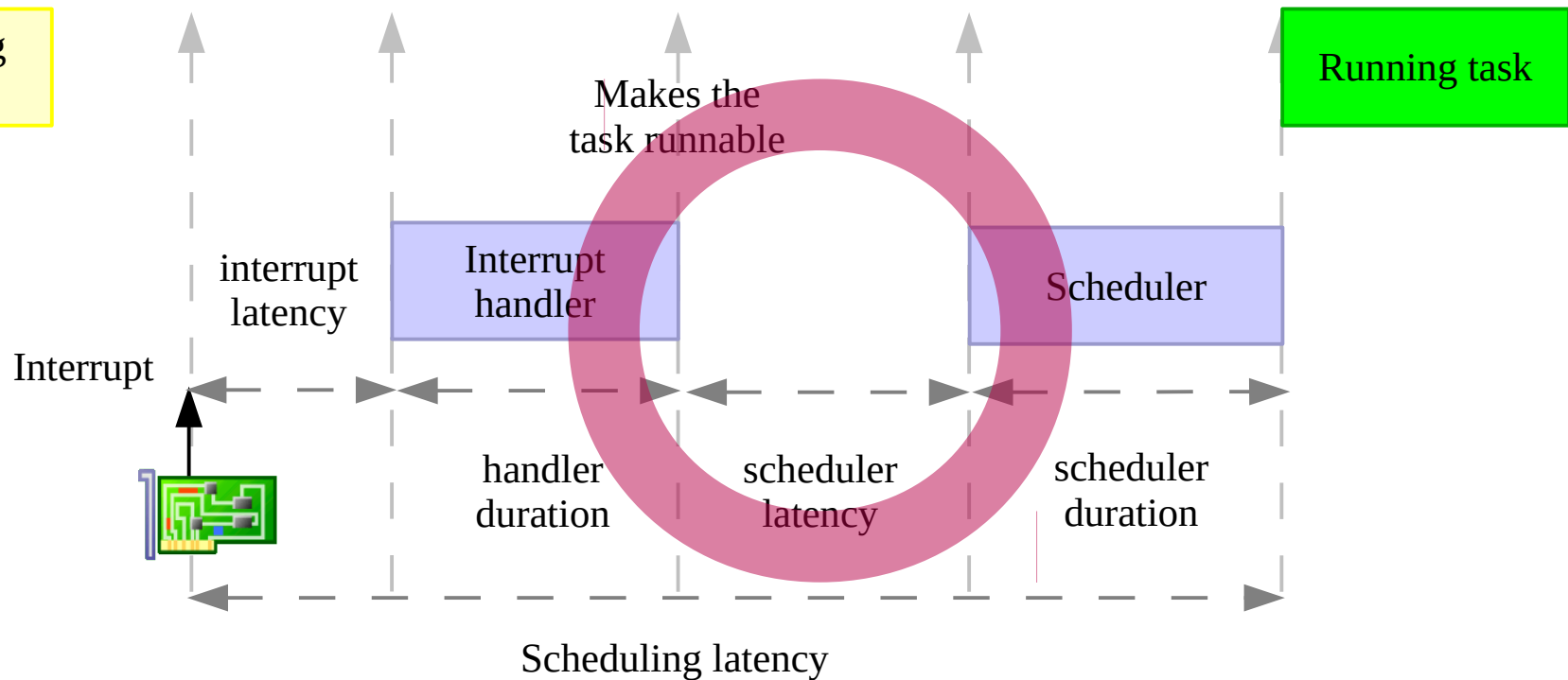
# Interrupt handler implementation

- ▶ In Linux, many interrupt handlers are split in two parts
  - ▶ A top-half, started by the CPU as soon as interrupt are enabled. It runs with the interrupt line disabled and is supposed to complete as quickly as possible.
  - ▶ A bottom-half, scheduled by the top-half, which starts after all pending top-half have completed their execution.
- ▶ Therefore, for real-time critical interrupts, bottom-half shouldn't be used: their execution is delayed by all other interrupts in the system.

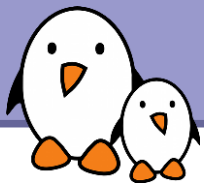




# Scheduler latency

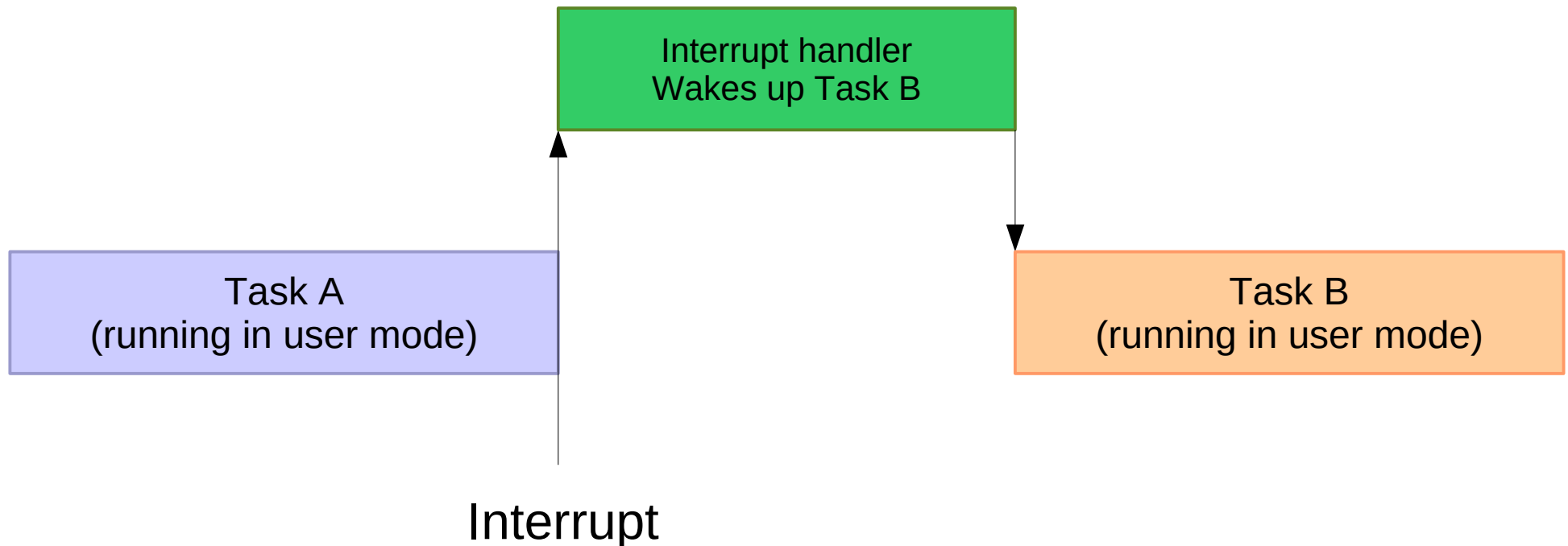


Time elapsed before executing the scheduler



# Understanding preemption (1)

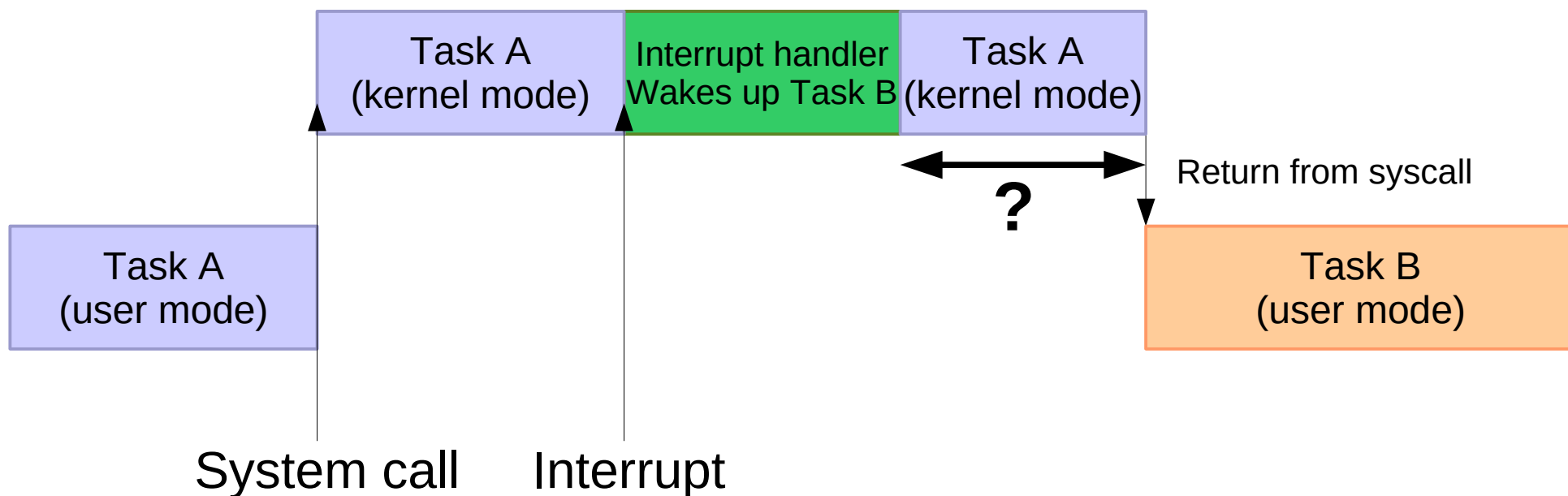
- ▶ The Linux kernel is a preemptive operating system
- ▶ When a task runs in user-space mode and gets interrupted by an interrupt, if the interrupt handler wakes up another task, this task can be scheduled as soon as we return from the interrupt handler.

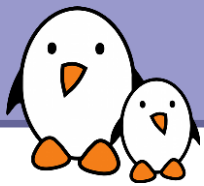




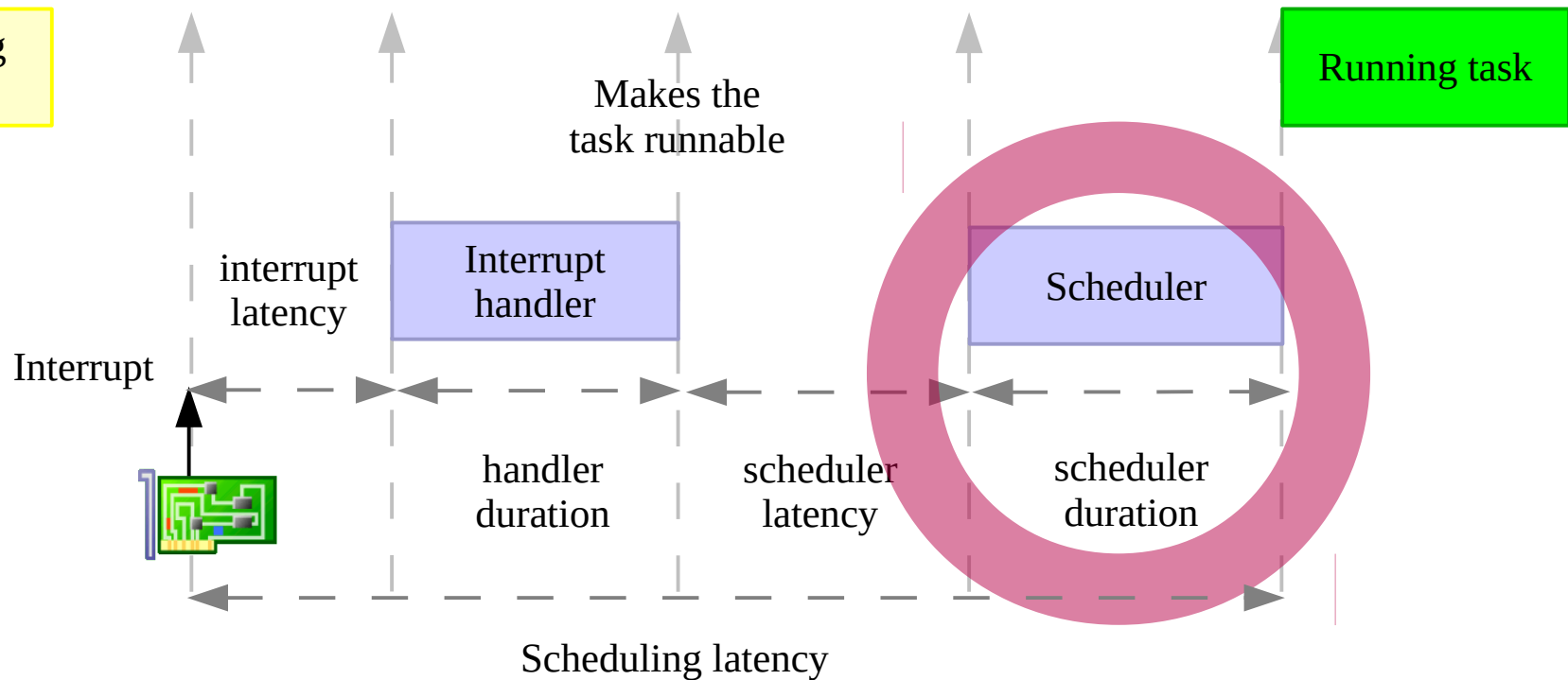
# Understanding preemption (2)

- ▶ However, when the interrupt comes while the task is executing a system call, this system call has to finish before another task can be scheduled.
- ▶ By *default*, the Linux kernel does not do *kernel preemption*.
- ▶ This means that the time before which the scheduler will be called to schedule another task is unbounded.



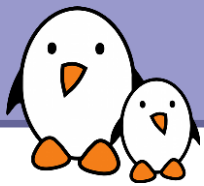


# Scheduler duration



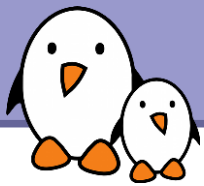
Time taken to execute the scheduler and switch to the new task.



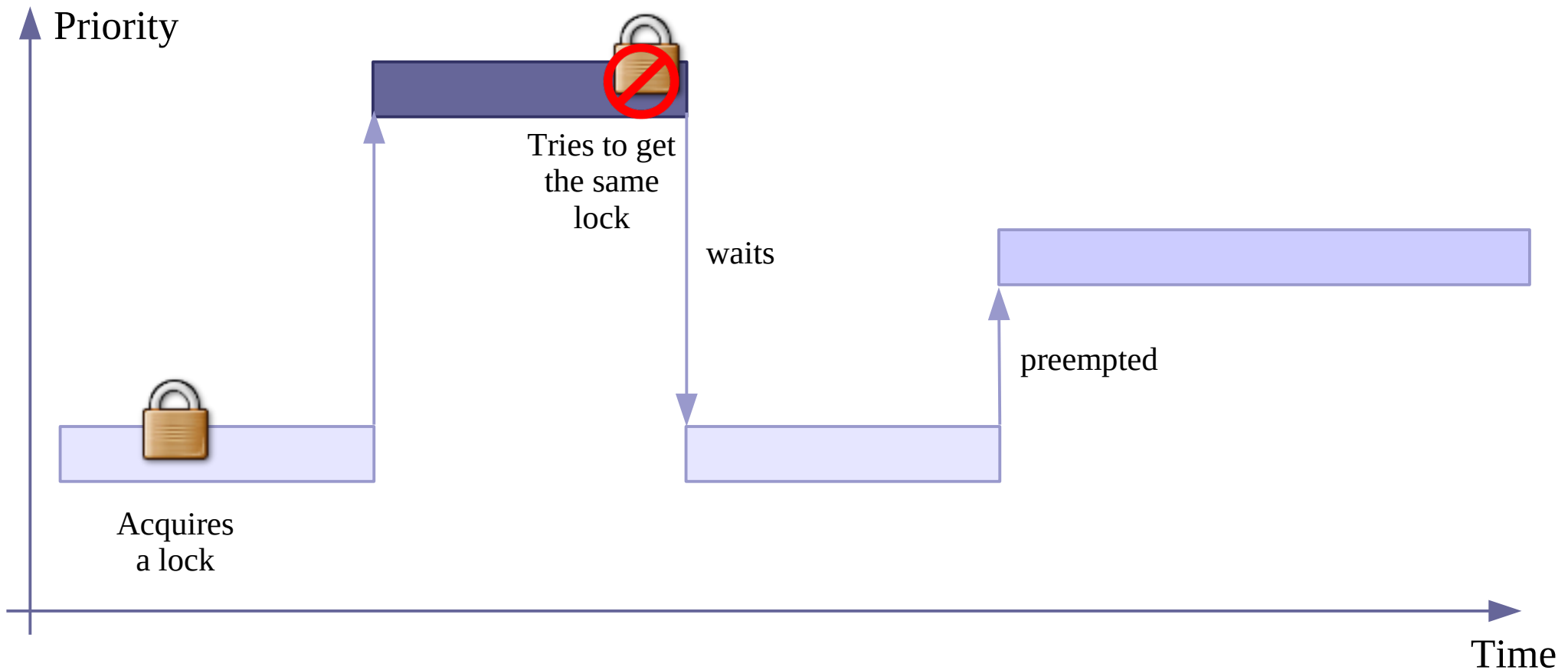


# Other non-deterministic mechanisms

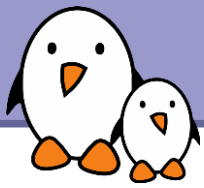
- ▶ Outside of the critical path detailed previously, other non-deterministic mechanisms of Linux can affect the execution time of real-time tasks
- ▶ Linux is highly based on virtual memory, as provided by an MMU, so that memory is allocated on demand. Whenever an application accesses code or data for the first time, it is loaded on demand, which can create huge delays.
- ▶ Many C library services or kernel services are not designed with real-time constraints in mind.



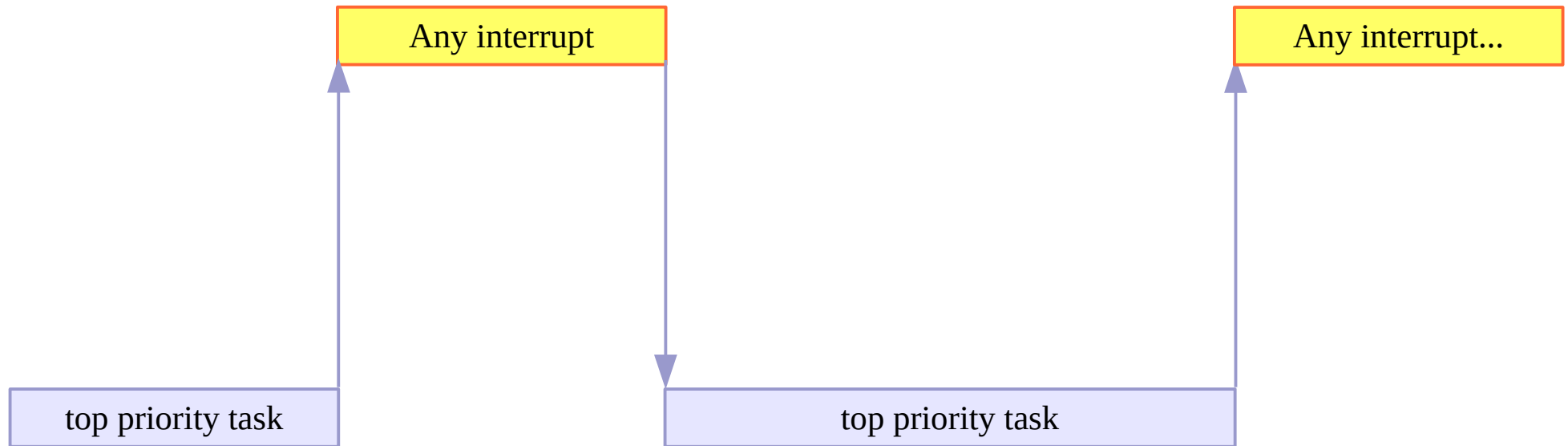
# Issue: priority inversion



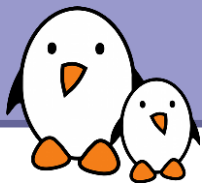
A process with more priority can preempt a process holding the lock. The top priority process could wait for a very long time.



# Issue: interrupt inversion

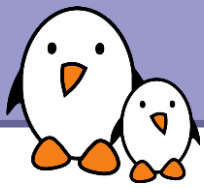


Even your top priority task can be “preempted” by any interrupt handler, even for interrupts feeding tasks with lower priority.



# The PREEMPT\_RT project

- ▶ Long-term project lead by Linux kernel developers Ingo Molnar, Thomas Gleixner and Steven Rostedt
  - ▶ <https://rt.wiki.kernel.org>
- ▶ The goal is to gradually improve the Linux kernel regarding real-time requirements and to get these improvements merged into the mainline kernel
  - ▶ PREEMPT\_RT development works very closely with the mainline development
- ▶ Many of the improvements designed, developed and debugged inside PREEMPT\_RT over the years are now part of the mainline Linux kernel
  - ▶ The project is a long-term branch of the Linux kernel that ultimately should disappear as everything will have been merged



# Improvements in the mainline kernel

- ▶ Coming from the PREEMPT\_RT project
- ▶ Since the beginning of 2.6
  - ▶ O(1) scheduler
  - ▶ Kernel preemption
  - ▶ Better POSIX real-time API support
- ▶ Since 2.6.18
  - ▶ Priority inheritance support for mutexes
- ▶ Since 2.6.21
  - ▶ High-resolution timers
- ▶ Since 2.6.30
  - ▶ Threaded interrupts
- ▶ Since 2.6.33
  - ▶ Spinlock annotations

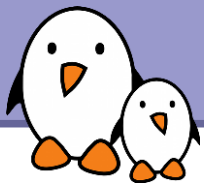


# New preemption options in Linux 2.6

2 new preemption models offered by standard Linux 2.6:

## Preemption Model

- |  |                   |
|--|-------------------|
| ○ No Forced Preemption (Server)            | PREEMPT_NONE      |
| ● Voluntary Kernel Preemption (Desktop)    | PREEMPT_VOLUNTARY |
| ○ Preemptible Kernel (Low-Latency Desktop) | PREEMPT           |

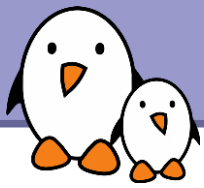


# 1<sup>st</sup> option: no forced preemption

## `CONFIG_PREEMPT_NONE`

Kernel code (interrupts, exceptions, system calls) never preempted.  
Default behavior in standard kernels.

- ▶ Best for systems making intense computations, on which overall throughput is key.
- ▶ Best to reduce task switching to maximize CPU and cache usage (by reducing context switching).
- ▶ Still benefits from some Linux 2.6 improvements: O(1) scheduler, increased multiprocessor safety (work on RT preemption was useful to identify hard to find SMP bugs).
- ▶ Can also benefit from a lower timer frequency (100 Hz instead of 250 or 1000).



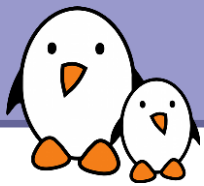
## 2<sup>nd</sup> option: voluntary kernel preemption

### CONFIG\_PREEMPT\_VOLUNTARY

Kernel code can preempt itself

- ▶ Typically for desktop systems, for quicker application reaction to user input.
- ▶ Adds explicit rescheduling points throughout kernel code.
- ▶ Minor impact on throughput.



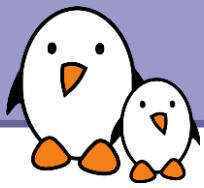


## 3<sup>rd</sup> option: preemptible kernel

### CONFIG\_PREEMPT

Most kernel code can be involuntarily preempted at any time. When a process becomes runnable, no more need to wait for kernel code (typically a system call) to return before running the scheduler.

- ▶ Exception: kernel critical sections (holding spinlocks), but a rescheduling point occurs when exiting the outer critical section, in case a preemption opportunity would have been signaled while in the critical section.
- ▶ Typically for desktop or embedded systems with latency requirements in the milliseconds range.
- ▶ Still a relatively minor impact on throughput.



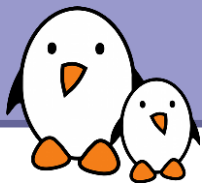
# Priority inheritance

- ▶ One classical solution to the priority inversion problem is called priority inheritance
  - ▶ The idea is that when a task of a low priority holds a lock requested by an higher priority task, the priority of the first task gets temporarily raised to the priority of the second task : it has *inherited* its priority.
- ▶ In Linux, since 2.6.18, mutexes support priority inheritance
- ▶ In userspace, priority inheritance must be explicitly enabled on a per-mutex basis.



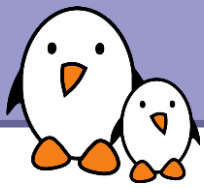
# High resolution timers

- ▶ The resolution of the timers used to be bound to the resolution of the regular system tick
  - ▶ Usually 100 Hz or 250 Hz, depending on the architecture and the configuration
  - ▶ A resolution of only 10 ms or 4 ms.
  - ▶ Increasing the regular system tick frequency is not an option as it would consume too much resources
- ▶ The high-resolution timers infrastructure, merged in 2.6.21, allows to use the available hardware timers to program interrupts at the right moment.
  - ▶ Hardware timers are multiplexed, so that a single hardware timer is sufficient to handle a large number of software-programmed timers.
  - ▶ Usable directly from user-space using the usual timer APIs



# Threaded interrupts

- ▶ To solve the interrupt inversion problem, PREEMPT\_RT has introduced the concept of threaded interrupts
- ▶ The interrupt handlers run in normal kernel threads, so that the priorities of the different interrupt handlers can be configured
- ▶ The real interrupt handler, as executed by the CPU, is only in charge of masking the interrupt and waking-up the corresponding thread
- ▶ The idea of threaded interrupts also allows to use sleeping spinlocks (see later)
- ▶ Merged since 2.6.30, the conversion of interrupt handlers to threaded interrupts is not automatic : drivers must be modified
- ▶ In PREEMPT\_RT, all interrupt handlers are switched to threaded interrupts



# PREEMPT\_RT specifics



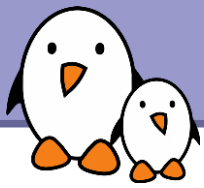
# CONFIG\_PREEMPT\_RT (1)

- ▶ The PREEMPT\_RT patch adds a new « level » of preemption, called CONFIG\_PREEMPT\_RT
- ▶ This level of preemption replaces all kernel spinlocks by mutexes (or so-called sleeping spinlocks)
  - ▶ Instead of providing mutual exclusion by disabling interrupts and preemption, they are just normal locks : when contention happens, the process is blocked and another one is selected by the scheduler
  - ▶ Works well with threaded interrupts, since threads can block, while usual interrupt handlers could not
  - ▶ Some core, carefully controlled, kernel spinlocks remain as normal spinlocks



# CONFIG\_PREEMPT\_RT (2)

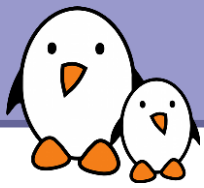
- ▶ With CONFIG\_PREEMPT\_RT, virtually all kernel code becomes preemptible
  - ▶ An interrupt can occur at any time, when returning from the interrupt handler, the woken up process can start immediately
- ▶ This is the last big part of PREEMPT\_RT that isn't fully in the mainline kernel yet
  - ▶ Part of it has been merged in 2.6.33 : the spinlock annotations. The spinlocks that must remain as spinning spinlocks are now differentiated from spinlocks that can be converted to sleeping spinlocks. This has reduced a lot the PREEMPT\_RT patch size !



# Threaded interrupts

- ▶ The mechanism of threaded interrupts in PREEMPT\_RT is still different from the one merged in mainline
- ▶ In PREEMPT\_RT, all interrupt handlers are unconditionally converted to threaded interrupts.
- ▶ This is a temporary solution, until interesting drivers in mainline get gradually converted to the new threaded interrupt API that has been merged in 2.6.30.





# Setting up PREEMPT\_RT



# PREEMPT\_RT setup (1)

- ▶ PREEMPT\_RT is delivered as a patch against the mainline kernel
  - ▶ Best to have a board supported by the mainline kernel, otherwise the PREEMPT\_RT patch may not apply and may require some adaptations
- ▶ Many official kernel releases are supported, but not all. For example, 2.6.31 and 2.6.33 are supported, but not 2.6.32.
- ▶ Quick set up
  - ▶ Download and extract mainline kernel
  - ▶ Download the corresponding PREEMPT\_RT patch
  - ▶ Apply it to the mainline kernel tree



# PREEMPT\_RT setup (2)

- ▶ In the kernel configuration, be sure to enable
  - ▶ `CONFIG_PREEMPT_RT`
  - ▶ High-resolution timers
- ▶ Compile your kernel, and boot
- ▶ You are now running the real-time Linux kernel
- ▶ Of course, some system configuration remains to be done, in particular setting appropriate priorities to the interrupt threads, which depend on your application.



# Real-time application development



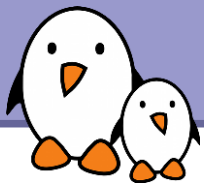
# Development and compilation

- ▶ No special library is needed, the POSIX realtime API is part of the standard C library
- ▶ The glibc or eglibc C libraries are recommended, as the support of some real-time features is not available yet in uClibc
  - ▶ Priority inheritance mutexes or NPTL on some architectures, for example
- ▶ Compile a program
  - ▶ `ARCH-linux-gcc -o myprog myprog.c -lrt`
- ▶ To get the documentation of the POSIX API
  - ▶ Install the `manpages-posix-dev` package
  - ▶ Run `man functionname`



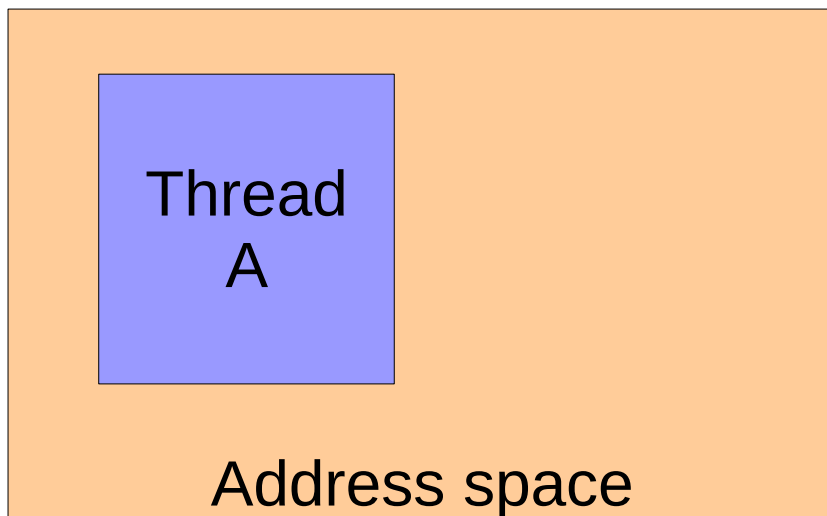
# Process, thread ?

- ▶ Confusion about the terms «process», «thread» and «task»
- ▶ In Unix, a process is created using `fork()` and is composed of
  - ▶ An address space, which contains the program code, data, stack, shared libraries, etc.
  - ▶ One thread, that starts executing the `main()` function.
  - ▶ Upon creation, a process contains one thread
- ▶ Additional threads can be created inside an existing process, using `pthread_create()`
  - ▶ They run in the same address space as the initial thread of the process
  - ▶ They start executing a function passed as argument to `pthread_create()`

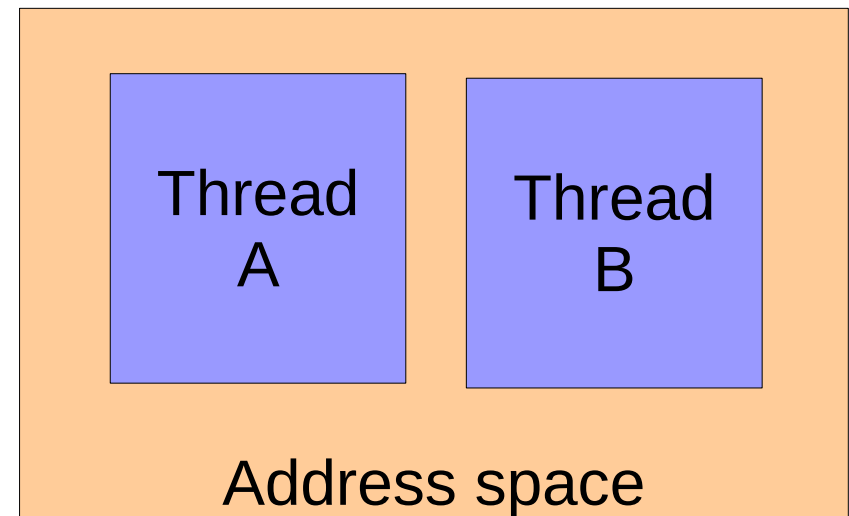


# Process, thread: kernel point of view

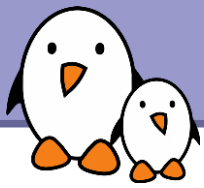
- ▶ The kernel represents each thread running in the system by a structure of type `task_struct`
- ▶ From a scheduling point of view, it makes no difference between the initial thread of a process and all additional threads created dynamically using `pthread_create()`



Process after `fork()`



Same process after `pthread_create()`



# Creating threads

- ▶ Linux support the POSIX thread API
- ▶ To create a new thread
  - ▶ **pthread\_create**(pthread\_t \*thread,  
pthread\_attr\_t \*attr,  
void \*(\*routine)(\*void\*),  
void \*arg);
  - ▶ The new thread will run in the same address space, but will be scheduled independently
- ▶ Exiting from a thread
  - ▶ **pthread\_exit**(void \*value\_ptr);
- ▶ Waiting for a thread termination
  - ▶ **pthread\_join**(pthread\_t \*thread, void \*\*value\_ptr);





# Scheduling classes (1)

- ▶ The Linux kernel scheduler support different scheduling classes
- ▶ The default class, in which processes are started by default is a time-sharing class
  - ▶ All processes, regardless of their priority, get some CPU time
  - ▶ The proportion of CPU time they get is dynamic and affected by the nice value, which ranges from -20 (highest) to 19 (lowest). Can be set using the nice or renice commands
- ▶ The real-time classes **SCHED\_FIFO** and **SCHED\_RR**
  - ▶ The highest priority process gets all the CPU time, until it blocks.
  - ▶ In SCHED\_RR, round-robin scheduling between the processes of the same priority. All must block before lower priority processes get CPU time.
  - ▶ Priorities ranging from 0 (lowest) to 99 (highest)



# Scheduling classes (2)

- ▶ Before creating a thread :

```
struct sched_param parm;  
pthread_attr_t attr;
```

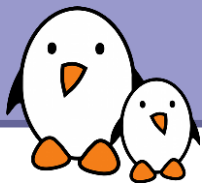
```
pthread_attr_init(&attr);  
pthread_attr_setinheritsched(&attr,  
                             PTHREAD_EXPLICIT_SCHED);  
pthread_attr_setschedpolicy(&attr, SCHED_FIFO);  
parm.sched_priority = 42;  
pthread_attr_setschedparam(&attr, &parm);
```

- ▶ Then the thread can be created using `pthread_create()`, passing the `attr` structure.
- ▶ Several other attributes can be defined this way: stack size, etc.



# Memory locking

- ▶ In order to solve the non-determinism introduced by virtual memory, memory can be locked
  - ▶ Guarantee that the system will keep it allocated
  - ▶ Guarantee that the system has pre-loaded everything into memory
- ▶ `mlockall(MCL_CURRENT | MCL_FUTURE);`
  - ▶ Locks all the memory of the current address space, for currently mapped pages and pages mapped in the future
- ▶ Other, less useful parts of the API: `munlockall`, `mock`, `munlock`.
- ▶ Watch out for non-currently mapped pages
  - ▶ Stack pages
  - ▶ Dynamically-allocated memory



# Mutexes

- ▶ Allows mutual exclusion between two threads in the same address space
  - ▶ Initialization/destruction

```
pthread_mutex_init(pthread_mutex_t *mutex, const  
pthread_mutexattr_t *mutexattr);  
pthread_mutex_destroy(pthread_mutex_t *mutex);
```
  - ▶ Lock/unlock

```
pthread_mutex_lock(pthread_mutex_t *mutex);  
pthread_mutex_unlock(pthread_mutex_t *mutex);
```
- ▶ Priority inheritance must explicitly be activated

```
pthread_mutexattr_t attr;  
pthread_mutexattr_init (&attr);  
pthread_mutexattr_getprotocol  
(&attr, PTHREAD_PRIO_INHERIT);
```



# Timers

- ▶ **timer\_create**(clockid\_t clockid,  
                  struct sigevent \*evp,  
                  timer\_t \*timerid)
  - ▶ Create a timer. **clockid** is usually CLOCK\_MONOTONIC.  
**sigevent** defines what happens upon timer expiration : send a signal or start a function in a new thread. **timerid** is the returned timer identifier.
- ▶ **timer\_settime**(timer\_t timerid, int flags,  
                  struct itimerspec \*newvalue,  
                  struct itimerspec \*oldvalue)
  - ▶ Configures the timer for expiration at a given time.
- ▶ **timer\_delete**(timer\_t timerid), delete a timer
- ▶ **clock\_getres**( ), get the resolution of a clock
- ▶ Other functions: timer\_getoverrun(), timer\_gettime()



# Signals

- ▶ Signals are an asynchronous notification mechanism
- ▶ Notification occurs either
  - ▶ By the call of a signal handler. Be careful with the limitations of signal handlers!
  - ▶ By being unblocked from the **sigwait()**, **sigtimedwait()** or **sigwaitinfo()** functions. Usually better.
- ▶ Signal behaviour can be configured using **sigaction()**
- ▶ Mask of blocked signals can be changed with **pthread\_sigmask()**
- ▶ Delivery of a signal using **pthread\_kill()** or **tgkill()**
- ▶ All signals between **SIGRTMIN** and **SIGRTMAX**, 32 signals under Linux.



# Inter-process communication

## ▶ Semaphores

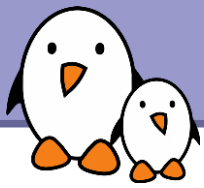
- ▶ Usable between different processes using named semaphores
- ▶ **sem\_open()**, **sem\_close()**, **sem\_unlink()**, **sem\_init()**, **sem\_destroy()**, **sem\_wait()**, **sem\_post()**, etc.

## ▶ Message queues

- ▶ Allows processes to exchange data in the form of messages.
- ▶ **mq\_open()**, **mq\_close()**, **mq\_unlink()**, **mq\_send()**, **mq\_receive()**, etc.

## ▶ Shared memory

- ▶ Allows processes to communicate by sharing a segment of memory
- ▶ **shm\_open()**, **ftruncate()**, **mmap()**, **munmap()**, **close()**, **shm\_unlink()**



# Debugging real-time latencies





# ftrace - Kernel function tracer

New infrastructure that can be used for debugging or analyzing latencies and performance issues in the kernel.

- ▶ Developed by Steven Rostedt. Merged in 2.6.27.  
For earlier kernels, can be found from the rt-preempt patches.
- ▶ Very well documented in [Documentation/ftrace.txt](#)
- ▶ Negligible overhead when tracing is not enabled at run-time.
- ▶ Can be used to trace any kernel function!
- ▶ See our video of Steven's tutorial at OLS 2008:  
<https://bootlin.com/community/videos/conferences/>



# Using ftrace

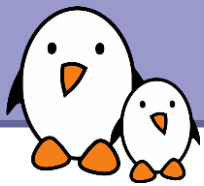
- ▶ Tracing information available through the debugfs virtual fs (`CONFIG_DEBUG_FS` in the `Kernel Hacking` section)
- ▶ Mount this filesystem as follows:  
`mount -t debugfs nodev /debug`
- ▶ When tracing is enabled (see the next slides), tracing information is available in `/debug/tracing`.
- ▶ Check available tracers  
in `/debug/tracing/available_tracers`



# Scheduling latency tracer

## CONFIG\_SCHED\_TRACER (Kernel Hacking section)

- ▶ Maximum recorded time between waking up a top priority task and its scheduling on a CPU, expressed in  $\mu$ s.
- ▶ Check that `wakeup` is listed in `/debug/tracing/available_tracers`
- ▶ To select, reset and enable this tracer:  
`echo wakeup > /debug/tracing/current_tracer`  
`echo 0 > /debug/tracing/tracing_max_latency`  
`echo 1 > /debug/tracing/tracing_enabled`
- ▶ Let your system run, in particular real-time tasks.  
Example: `chrt -f 5 sleep 1`
- ▶ Disable tracing:  
`echo 0 > /debug/tracing/tracing_enabled`
- ▶ Read the maximum recorded latency and the corresponding trace:  
`cat /debug/tracing/tracing_max_latency`



# Xenomai



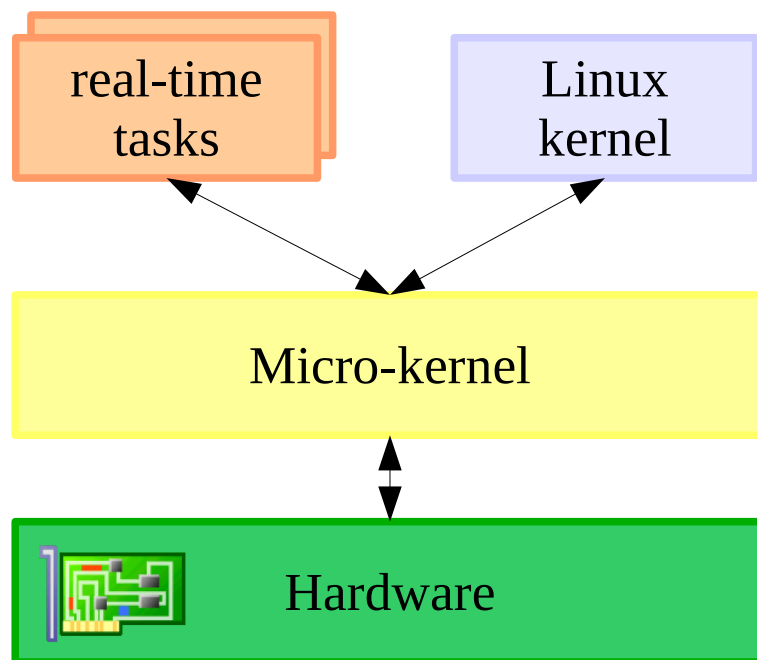
# Linux hard real-time extensions

## Three generations

- ▶ RTLinux
- ▶ RTAI
- ▶ Xenomai

## A common principle

- ▶ Add an extra layer between the hardware and the Linux kernel, to manage real-time tasks separately.





# Xenomai project

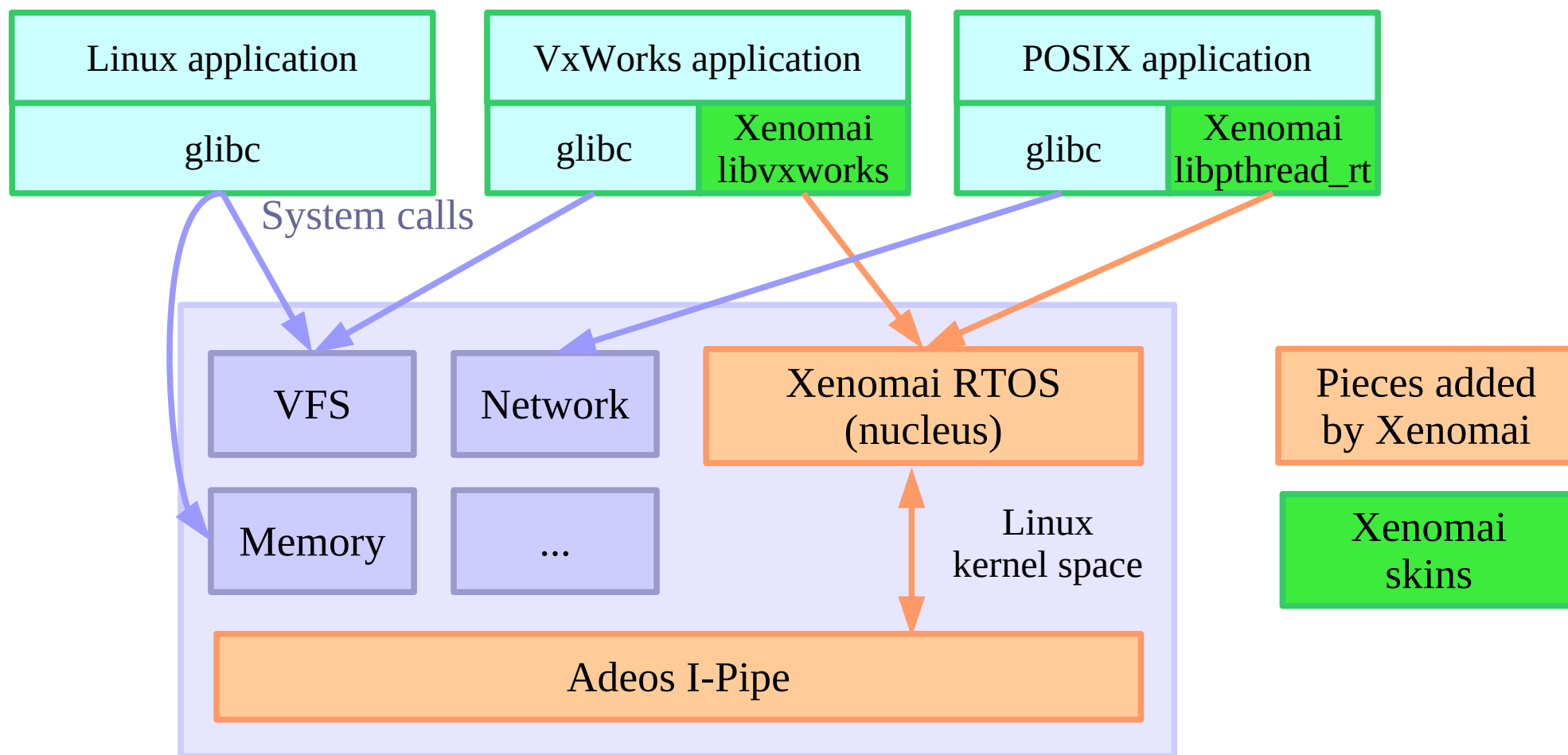


<http://www.xenomai.org/>

- ▶ Started in 2001 as a project aiming at emulating traditional RTOS.
- ▶ Initial goals: facilitate the porting of programs to GNU / Linux.
- ▶ Initially related to the RTAI project (as the RTAI / fusion branch), now independent.
- ▶ Skins mimicking the APIs of traditional RTOS such as VxWorks, pSOS+, and VRTXsa as well as the POSIX API, and a “native” API.
- ▶ Aims at working both as a co-kernel and on top of PREEMPT\_RT in the upcoming 3.0 branch.
- ▶ Will never be merged in the mainline kernel.



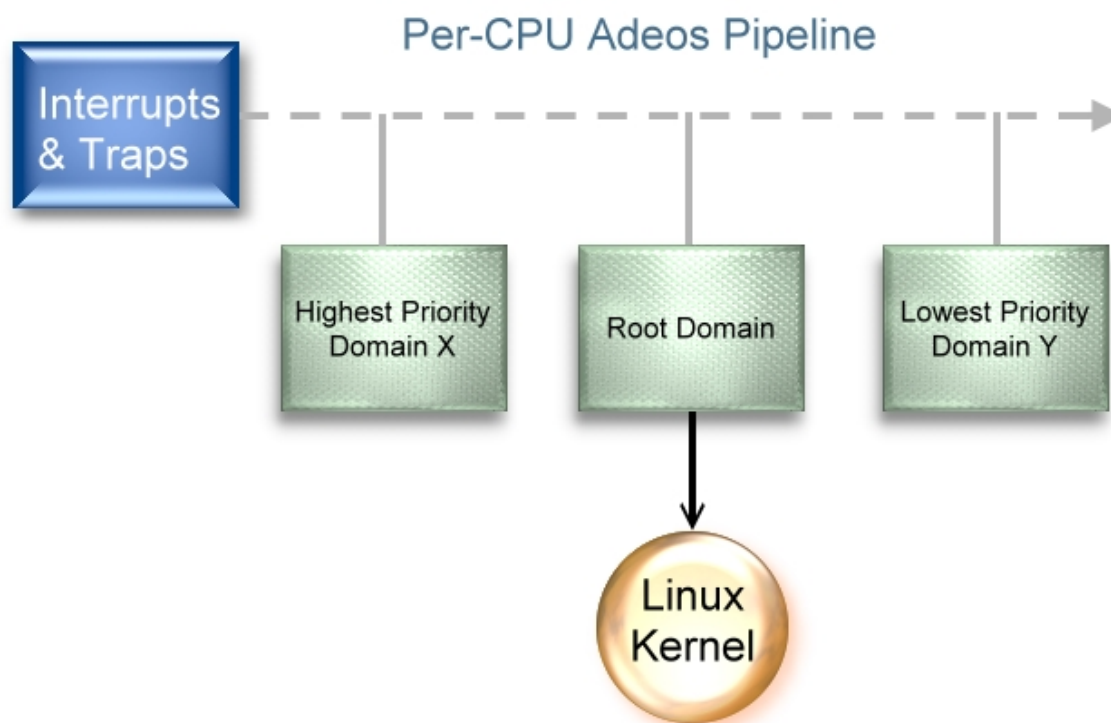
# Xenomai architecture





# The Adeos interrupt pipeline abstraction

- ▶ From Adeos point of view, guest OSES are prioritized domains.
- ▶ For each event (interrupts, exceptions, syscalls, etc...), the various domains may handle the event or pass it down the pipeline.

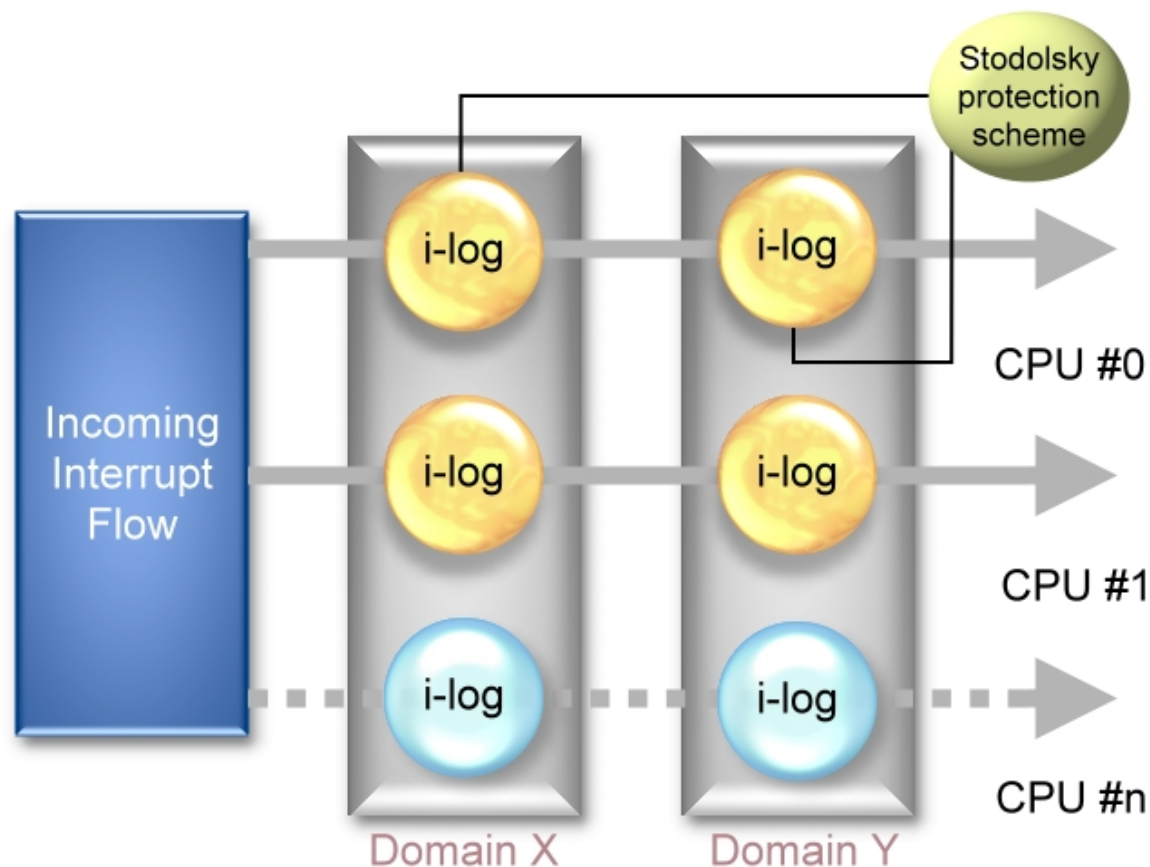






# Adeos virtualized interrupts disabling

- ▶ Each domain may be “stalled”, meaning that it does not accept interrupts.
- ▶ Hardware interrupts are not disabled however (except for the domain leading the pipeline), instead the interrupts received during that time are logged and replayed when the domain is unstalled.





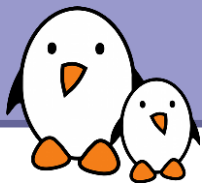
# Adeos additional features

- ▶ The Adeos I-pipe patch implement additional features, essential for the implementation of the Xenomai real-time extension:
  - ▶ Disables on-demand mapping of kernel-space vmalloc/ioremap areas.
  - ▶ Disables copy-on-write when real-time processes are forking.
  - ▶ Allow subscribing to event allowing to follow progress of the Linux kernel, such as Linux system calls, context switches, process destructions, POSIX signals, FPU faults.
  - ▶ On the ARM architectures, integrates the FCSE patch, which allows to reduce the latency induced by cache flushes during context switches.



# Xenomai features

- ▶ Factored real-time core with skins implementing various real-time APIs
- ▶ Seamless support for hard real-time in user-space
- ▶ No second-class citizen, all ports are equivalent feature-wise
- ▶ Xenomai support is as much as possible independent from the Linux kernel version (backward and forward compatible when reasonable)
- ▶ Each Xenomai branch has a stable user/kernel ABI
- ▶ Timer system based on hardware high-resolution timers
- ▶ Per-skin time base which may be periodic
- ▶ RTDM skin allowing to write real-time drivers



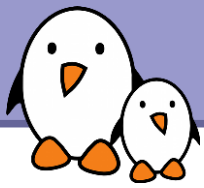
# Xenomai user-space real-time support.

- ▶ Xenomai supports real-time in user-space on 5 architectures, including 32 and 64 bits variants.
- ▶ Two modes are defined for a thread
  - ▶ the primary mode, where the thread is handled by Xenomai scheduler
  - ▶ the secondary mode, when it is handled by Linux scheduler.
- ▶ Thanks to the services of the Adeos I-pipe service, Xenomai system calls are defined.
- ▶ A thread migrates from secondary mode to primary mode when such a system call is issued
- ▶ It migrates from primary mode to secondary mode when a Linux system call is issued, or to handle gracefully exceptional events such as exceptions or Linux signals.



# Real Time Driver Model (RTDM)

- ▶ An approach to unify the interfaces for developing device drivers and associated applications under real-time Linux
  - ▶ An API very similar to the native Linux kernel driver API
- ▶ Allows the development, in kernel space, of
  - ▶ Character-style device drivers
  - ▶ Network-style device drivers
- ▶ See the whitepaper on <http://www.xenomai.org/documentation/xenomai-2.4/pdf/RTDM-and-Applications.pdf>
- ▶ Current notable RTDM based drivers:
  - ▶ Serial port controllers;
  - ▶ RTnet UDP/IP stack;
  - ▶ RT socket CAN, drivers for CAN controllers;
  - ▶ Analogy, fork of the Comedy project, drivers for acquisition cards.



# Setting up Xenomai



# How to build Xenomai

- ▶ Download Xenomai sources at <http://download.gna.org/xenomai/stable/>
- ▶ Download one of the Linux versions supported by this release (see [ksrc/arch/<arch>/patches/](#))
- ▶ Since version 2.0, split kernel/user building model.
- ▶ Kernel uses a script called [script/prepare-kernel.sh](#) which integrates Xenomai kernel-space support in the Linux sources.
- ▶ Run the kernel configuration menu.



# Linux options for Xenomai configuration

File Edit Option Help

Option

- General setup
  - RCU Subsystem
  - ☐ Control Group support
  - ☐ Configure standard kernel fea
- ☒ Enable loadable module support
- Enable the block layer (NEW)
  - IO Schedulers
- Real-time sub-system
- System Type
  - Atmel AT91 System-on-Chip
- Bus support
  - ☐ PCCard (PCMCIA/CardBus) s
- Kernel Features
- Boot options
- CPU Power Management
- Floating point emulation
- Userspace binary formats
- Power management options

Option

- ☒ Xenomai
- ☒ Nucleus
  - ☒ Pervasive real-time support in user-space
  - ☒ Optimize as pipeline head
  - ☐ Extra scheduling classes
  - (32) Number of pipe devices
  - (512) Number of registry slots
  - (256) Size of the system heap (Kb)
  - (128) Size of the private stack pool (Kb)
  - (12) Size of private semaphores heap (Kb)
  - (12) Size of global semaphores heap (Kb)
  - ☒ Statistics collection

**Xenomai (XENOMAI)**

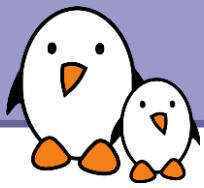
Xenomai is a real-time extension to the Linux kernel. Note that Xenomai relies on Adeos interrupt pipeline (CONFIG\_IPIPE option) to be enabled, so enabling this option selects the CONFIG\_IPIPE option.





# Xenomai user-space support

- ▶ User-space libraries are compiled using the traditional autotools
  - ▶ `./configure --target=arm-linux && make && make DESTDIR=/your/rootfs/ install`
- ▶ The `xeno-config` script, installed when installing Xenomai user-space support helps you compiling your own programs.
- ▶ See Xenomai's `examples` directory.
- ▶ Installation details may be found in the `README.INSTALL` guide.
- ▶ For an introduction on programming with the native API, see:  
<http://www.xenomai.org/documentation/branches/v2.3.x/pdf/Native-API-Tour-rev-C.pdf>
- ▶ For an introduction on programming with the POSIX API, see:  
[http://www.xenomai.org/index.php/Porting\\_POSIX\\_applications\\_to\\_Xenomai](http://www.xenomai.org/index.php/Porting_POSIX_applications_to_Xenomai)



# Developing applications on Xenomai



# The POSIX skin

- ▶ The POSIX skin allows to recompile without changes a traditional POSIX application so that instead of using Linux real-time services, it uses Xenomai services
  - ▶ Clocks and timers, condition variables, message queues, mutexes, semaphores, shared memory, signals, thread management
  - ▶ Good for existing code or programmers familiar with the POSIX API
- ▶ Of course, if the application uses any Linux service that isn't available in Xenomai, it will switch back to secondary mode
- ▶ To link an application against the POSIX skin

```
DESTDIR=/path/to/xenomai/  
export DESTDIR  
CFL=`$DESTDIR/bin/xeno-config --posix-cflags`  
LDF=`$DESTDIR/bin/xeno-config --posix-ldflags`  
ARCH-gcc $CFL -o rttest rttest.c $LDF
```



# Communication with a normal task

- ▶ If a Xenomai real-time application using the POSIX skin wishes to communicate with a separate non-real-time application, it must use the rtipc mechanism
- ▶ In the Xenomai application, create an IPCPROTO\_XDDP socket

```
socket(AF_RTIPC, SOCK_DGRAM, IPCPROTO_XDDP);
setsockopt(s, SOL_RTIPC, XDDP_SETLOCALPOOL, &poolsz,
sizeof(poolsz));
memset(&saddr, 0, sizeof(saddr));
saddr.sipc_family = AF_RTIPC;
saddr.sipc_port = MYAPPIDENTIFIER;
ret = bind(s, (struct sockaddr *)&saddr, sizeof(saddr));
```

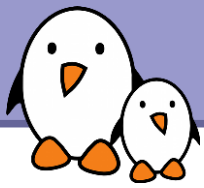
  - ▶ And then the normal socket API `sendto()` / `recvfrom()`
- ▶ In the Linux application
  - ▶ Open `/dev/rtpX`, where X is the XDDP port
  - ▶ Use `read()` and `write()`



# The native API (1)

- ▶ A Xenomai-specific API for developing real-time tasks
  - ▶ Usable both in user-space and kernel space. Development of tasks in user-space is the preferred way.
  - ▶ More coherent and more flexible API than the POSIX API. Easier to learn and understand. Certainly the way to go for new applications.
- ▶ Applications should include `<native/service.h>`, where service can be alarm, buffer, cond, event, heap, intr, misc, mutex, pipe, queue, sem, task, timer
- ▶ To compile applications :

```
DESTDIR=/path/to/xenomai/  
export DESTDIR  
CFL=`$DESTDIR/bin/xeno-config --xeno-cflags`  
LDF=`$DESTDIR/bin/xeno-config --xeno-ldflags`  
ARCH-gcc $CFL -o rttest rttest.c $LDF -lnative
```



# The native API (2)

## ▶ Task management services

- ▶ `rt_task_create()`, `rt_task_start()`,  
`rt_task_suspend()`, `rt_task_resume()`,  
`rt_task_delete()`, `rt_task_join()`, etc.

## ▶ Counting semaphore services

- ▶ `rt_sem_create()`, `rt_sem_delete()`, `rt_sem_p()`,  
`rt_sem_v()`, etc.

## ▶ Message queue services

- ▶ `rt_queue_create()`, `rt_queue_delete()`,  
`rt_queue_alloc()`, `rt_queue_free()`,  
`rt_queue_send()`, `rt_queue_receive()`, etc.

## ▶ Mutex services

- ▶ `rt_mutex_create()`, `rt_mutex_delete()`,  
`rt_mutex_acquire()`, `rt_mutex_release()`, etc.



# The native API (3)

## ▶ Alarm services

- ▶ `rt_alarm_create()`, `rt_alarm_delete()`,  
`rt_alarm_start()`, `rt_alarm_stop()`,  
`rt_alarm_wait()`, etc.

## ▶ Memory heap services

- ▶ Allows to share memory between processes and/or to pre-allocate a pool of memory

- ▶ `rt_heap_create()`, `rt_heap_delete()`,  
`rt_heap_alloc()`, `rt_heap_bind()`

## ▶ Condition variable services

- ▶ `rt_cond_create()`, `rt_cond_delete()`,  
`rt_cond_signal()`, `rt_cond_broadcast()`,  
`rt_cond_wait()`, etc.



# Xenomai and normal task communication

- ▶ Using *rt\_pipes*
- ▶ In the native Xenomai application, use the Pipe API
  - ▶ `rt_pipe_create()`, `rt_pipe_delete()`,  
`rt_pipe_receive()`, `rt_pipe_send()`,  
`rt_pipe_alloc()`, `rt_pipe_free()`
- ▶ In the normal Linux application
  - ▶ Open the corresponding `/dev/rtpX` file, the minor is specified at `rt_pipe_create()` time
  - ▶ Then, just `read()` and `write()` to the opened file

Xenomai application  
Uses the `rt_pipe_*`() API



Linux application  
`open("/dev/rtpX")`





# Questions ?

?

Thomas Petazzoni  
Bootlin, <https://bootlin.com>  
*[thomas.petazzoni@bootlin.com](mailto:thomas.petazzoni@bootlin.com)*



## Real-time - Testing PREEMT\_RT and Xenomai

### Using the Calao board

The board we are using is a Calao USB-A9263 device, which has an ARM AT91SAM9263 CPU clocked at 180 Mhz, 64 MB of RAM and 256 MB of NAND Flash.

The board is powered by the main USB connector, which is also used for the serial port console. It must also be connected to the development workstation using an Ethernet cable to the provided little USB-Ethernet adapter.

On your system, the following elements have already been configured :

- The terminal emulator *Minicom*, which allows to communicate with the board over the serial port. When the board is plugged-in, devices `/dev/ttyUSB0` and `/dev/ttyUSB1` are created. `/dev/ttyUSB0` is the JTAG port, that will not be used. `/dev/ttyUSB1` is the serial port. Minicom is already configured for `/dev/ttyUSB1`, speed 115200, no hardware flow control.
- A NFS server exports the directory `/home/tux/realtime/nfsroot/` to 192.168.42.2, which is the IP address of the board. This directory contains the root filesystem that will be used by the Linux system running on the ARM board.
- A TFTP server is installed, and exposes the contents of `/var/lib/tftpboot/`. This directory contains a *uImage* file, which is pre-compiled kernel image for the board.

Once the board is plugged-in, start Minicom. You will see the messages of the U-Boot bootloader, the messages of the Linux kernel, and finally the login prompt. The login is "root" and the password is empty.

### Testing the default kernel

In order to test the behaviour of the kernel, we have created a simple application in `/home/tux/realtime/nfsroot/root/rttest.c`, which arms a timer, and looks at what time the application is really woken up. It allows to measure the latency of the system.

The kernel we have compiled is a default 2.6.29 kernel, with high-resolution timers and no kernel preemption.

A CodeSourcery ARM toolchain has been pre-installed in `/home/tux/arm-2009q1`. So, let's add this toolchain to your PATH :

```
export PATH=/home/tux/arm-2009q1/bin:$PATH
```

Then, compile the example application :

```
arm-none-linux-gnueabi-gcc -o rttest rttest.c -lrt
```

Now, do the following tests:

- Test the program with nothing special and write down the results.
- Test your program and at the same time, add some workload



to the board, by running `dload 300` on the board, and using `netcat 192.168.0.100 5566` on your workstation when you see the message "Listening on any address 5566" in order to flood the network interface of the Calao board (where 192.168.0.100 is the IP address of the Calao board). A telnet server is running on the board, so you can get several shells on the board by running `telnet 192.168.42.2`. This allows you to run the `dload 300` command in one terminal, and test application in another.

- Test your program again with the workload, but by running the program in the `SCHED_FIFO` scheduling class at priority 80. To do so, you must first modify your program so that the testing takes place in a separate thread, created by `pthread_create()`. You'll have to use various functions of the `pthread` API to set up the thread in the correct scheduling class, at the correct priority.

## Testing the real-time patch

The sources of the mainline 2.6.29 kernel are available in

```
/home/tux/files/linux-2.6.29.6.tar.bz2
```

The `PREEMPT_RT` patch for 2.6.29.6 is available in

```
/home/tux/files/patch-2.6.29.6-rt24.bz2
```

Uncompress the kernel sources :

```
tar xjf /home/tux/files/linux-2.6.29.6.tar.bz2
```

Apply the `PREEMPT_RT` patch :

```
cd linux-2.6.29.6
```

```
bzcat /home/tux/files/patch-2.6.29.6-rt24.bz2 | patch -p1
```

Now, let's configure the kernel with a default configuration for the board :

```
make ARCH=arm usb-a9263_defconfig
```

Run the kernel configuration utility :

```
make ARCH=arm xconfig
```

And in the menus :

- Enable the `CONFIG_PREEMPT_RT` option. The kernel should now be fully preemptible, including in critical sections, and tasks can have higher priorities than interrupts.
- Enable the `CONFIG_HIGH_RES_TIMERS` option (in Kernel features), for high resolution timers.
- Enable the `CONFIG_ATMEL_TCLIB` option (in Device Drivers → Misc Devices), in order to have a clock source for the high resolution timers.

Compile the kernel:

```
make ARCH=arm CROSS_COMPILE=arm-none-linux-gnueabi- uImage
```

Finally, copy the kernel to the directory exported by the TFTP server :

```
sudo cp arch/arm/boot/uImage /var/lib/tftpboot/
```



Reboot your board and make sure that you are really running the new kernel.

Look at the output of the `ps` command : you now have threads for the interrupts handlers (threads named `[IRQ-XX]`) and threads for the bottom half handlers (threads named `[sirq-xxxx]`).

Repeat the tests with this real-time preemptible kernel and compare the results.

## Testing Xenomai scheduling latency

Xenomai is available in :

```
/home/tux/files/xenomai-2.5.2.tar.bz2
```

Extract the kernel sources again, in a separate directory (Xenomai and PREEMPT\_RT cannot be used at the same time).

Prepare the kernel for Xenomai compilation :

```
/home/tux/xenomai-2.5.2/scripts/prepare-kernel.sh -  
arch=arm --linux=/home/tux/linux-2.6.29.6
```

Start with a default kernel configuration for the Calao board :

```
make ARCH=ARM usb-a9263_defconfig
```

Then start the kernel configuration utility :

```
make ARCH=arm xconfig
```

And enable these options :

- `CONFIG_XENOMAI`
- `CONFIG_XENO_DRIVERS_TIMERBENCH`

Other options of interest (ARM specific) are:

- `CONFIG_ARM_FCSE_GUARANTEED`
- `CONFIG_XENO_HW_UNLOCKED_SWITCH`

Compile the kernel :

```
make ARCH=arm CROSS_COMPILE=arm-none-linux-gnueabi- uImage
```

And put it inside the directory exported by TFTP :

```
sudo cp arch/arm/boot/uImage /var/lib/tftpboot
```

Now boot the board with the new kernel.

We have already installed the user-space of Xenomai in `/home/tux/realtime/nfsroot/`. So, we can compile the same `rttest.c` application for Xenomai, using its POSIX skin :

```
DESTDIR=/home/tux/realtime/nfsroot/  
export DESTDIR  
CFL=`$DESTDIR/bin/xeno-config --posix-cflags`  
LDF=`$DESTDIR/bin/xeno-config --posix-ldflags`  
arm-none-linux-gnueabi-gcc $CFL -o rttest rttest.c $LDF
```

Re-run the tests, compare the results.

## Using the native Xenomai API

Now, we will start using the native Xenomai API. Your goal is to create a program that starts a Xenomai task, using the task



management Xenomai service. This task will loop infinitely, but will use the Xenomai alarm API to block and be woken up at regular intervals of 1 second. For the moment, every time the task is woken up, just print a message using the `printf()` function. Note that this is not correct from a real-time perspective: `printf()` is a function that uses Linux services, therefore our task will switch in secondary mode every time we call this function.

To compile your application :

```
DESTDIR=/home/tux/realtime/nfsroot/  
export DESTDIR  
CFL=`$DESTDIR/bin/xeno-config --xeno-cflags`  
LDF=`$DESTDIR/bin/xeno-config --xeno-ldflags`  
arm-none-linux-gnueabi-gcc $CFL -o xenotest xenotest.c\  
-lnative $LDF
```

Once this application works, we will improve it to demonstrate the usage of pipes to communicate with a non-Xenomai application. To do so, first improve your Xenomai application so that it creates a pipe. Then, at each iteration of the waiting loop, increment a counter, and using the `rt_pipe_write()` function, write the current value of the counter into the pipe.

Write another C application that opens `/dev/rtpX`, and infinitely reads integers from this device and display them on the screen.