

XtratuM Hypervisor for LEON3

Volume 2: User Manual

Miguel Masmano, Ismael Ripoll, Alfons Crespo

February, 2011

Reference: xm-3-usermanual-022c

This page is intentionally left blank.

DOCUMENT CONTROL PAGE

TITLE: XtratuM Hypervisor for LEON3 : Volume 2: User Manual

AUTHOR/S: Miguel Masmano, Ismael Ripoll, Alfons Crespo

LAST PAGE NUMBER: 119

VERSION OF SOURCE CODE: XtratuM 3 for LEON3 ()

REFERENCE ID: xm-3-usermanual-022c

SUMMARY: This guide describes the fundamental concepts of the XtratuM hypervisor, its API.

DISCLAIMER: This documentation is currently under active development. Therefore, no explicit or implied warranties in respect of any properties, including, but not limited to, correctness and fitness for purpose. Contributions of material, suggestions and corrections are welcome.

REFERENCING THIS DOCUMENT:

```
@techreport {xm-3-usermanual-022c,
  title = {XtratuM Hypervisor for LEON3 : Volume 2: User Manual},
  authors = {Miguel Masmano & Ismael Ripoll & Alfons Crespo},
  institution = {Universidad polit cnica de Valencia},
  number = {xm-3-usermanual-022c},
  date={February, 2011},
}
```

Copyright   February, 2011 Miguel Masmano, Ismael Ripoll and Alfons Crespo

Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.2 or any later version published by the Free Software Foundation; with no Invariant Sections, no Front-Cover Texts, and no Back-Cover Texts. A copy of the license is included in the section entitled "GNU Free Documentation License".

Changes:

Version	Date	Comments
0.1	February, 2010	[xm-3-usermanual-022] Initial document, based on the document XM-usermanual-002j. Code version 3.1.0. <ul style="list-style-type: none"> • Major changes in the ABI chapter (Binary interfaces). • Major changes in the tools chapter. • Internal object programming reworked. No impact on the libxm API.
0.2	May, 2010	[xm-3-usermanual-022b]. Code version 3.1.2 <ul style="list-style-type: none"> • Multi-plan support. Sections 2.5.1 and 5.7. • The PIT (Partition Information Table) has been removed. The PIT fields have been moved to the PCT. The PCT is not read-only. • Message notification support removed.
0.3	October, 2010	[????]. <ul style="list-style-type: none"> • Global line numbering added.

Version	Date	Comments
0.4	February, 2011	[xm-3-usermanual-022c]. Code version 3.2.0 •

Several typos corrected.

Contents

Preface	ix
1 Introduction	1
1.1 History	2
2 XtratuM Architecture	5
2.1 System operation	6
2.2 Partition operation	7
2.3 System partitions	8
2.4 Names and identifiers	8
2.5 Partition scheduling	9
2.5.1 Multiple scheduling plans	12
2.6 Inter-partition communications (IPC)	13
2.7 Health monitor (HM)	14
2.7.1 HM Events	15
2.7.2 HM Actions	16
2.7.3 HM Configuration	17
2.7.4 HM notification	17
2.8 Access to devices	17
2.9 Traps, interrupts and exceptions	18
2.9.1 Traps	18
2.9.2 Interrupts	19
2.10 Traces	19
2.11 Clocks and timers	20
2.12 Status	21
2.13 Summary	22
3 Developing Process Overview	23
3.1 Development at a glance	24
3.2 Building XtratuM	25

3.3	System configuration	26
3.4	Compiling partition code	27
3.5	Passing parameters to the partitions: customisation files	28
3.6	Building the final system image	28
4	Building XtratuM	29
4.1	Developing environment	29
4.2	Compile XtratuM Hypervisor	29
4.3	Generating binary a distribution	31
4.4	Installing a binary distribution	31
4.5	Compile the Hello World! partition	33
4.6	XtratuM directory tree	34
5	Partition Programming	35
5.1	Implementation requirements	35
5.2	XAL development environment	36
5.3	Partition definition	38
5.4	The “Hello World” example	39
5.4.1	Included headers	44
5.5	Partition reset	45
5.6	System reset	45
5.7	Scheduling	45
5.7.1	Slot identification	45
5.7.2	Managing scheduling plans	46
5.8	Console output	46
5.9	Inter-partition communication	47
5.9.1	Message notification	48
5.10	Peripheral programming	48
5.11	Traps, interrupts and exceptions	49
5.11.1	Traps	49
5.11.2	Interrupts	49
5.11.3	Exceptions	51
5.12	Clock and timer services	52
5.12.1	Execution time clock	52
5.13	Processor management	52
5.13.1	Managing stack context	53
5.14	Tracing	53
5.14.1	Trace messages	53

5.14.2	Reading traces	55
5.14.3	Configuration	55
5.15	System and partition status	55
5.16	Memory management	56
5.17	Releasing the processor	57
5.18	Partition customisation files	57
5.19	Assembly programming	58
5.19.1	The object interface	59
5.20	Manpages summary	59
6	Binary Interfaces	63
6.1	Data representation	63
6.2	Hypercall mechanism	64
6.3	Executable formats overview	64
6.4	Partition ELF format	65
6.4.1	Partition image header	65
6.4.2	Partition control table (PCT)	67
6.5	XEF format	68
6.5.1	Compression algorithm	71
6.6	Container format	71
7	Booting	75
7.1	Boot configuration	76
8	Configuration	79
8.1	XtratuM source code configuration (menuconfig)	79
8.2	Resident software source code configuration (menuconfig)	81
8.2.1	Memory requirements	82
8.3	Hypervisor configuration file (XM_CF)	83
8.3.1	Data representation and XPath syntax	83
8.3.2	The root element: /SystemDescription	85
8.3.3	The /SystemDescription/XMHypervisor element	85
8.3.4	The /SystemDescription/HwDescription element	87
8.3.5	The /SystemDescription/ResidentSw element	88
8.3.6	The /SystemDescription/PartitionTable/Partition element	88
8.3.7	The /SystemDescription/Channels element	90
9	Tools	91
9.1	XML configuration parser (xmcparser)	91

9.1.1	xmcparser	92
9.2	ELF to XEF (xmeformat)	92
9.2.1	xmeformat	92
9.3	Container builder (xmpack)	94
9.3.1	xmpack	94
9.4	Bootable image creator (rswbuild)	96
9.4.1	rswbuild	96
10	Security issues	97
10.1	Invoking a hypercall from libXM	97
10.2	Preventing covert/side channels due to scheduling slot overrun	97
A	XML Schema Definition	99
A.1	XML Schema file	99
A.2	Configuration file example	107
	GNU Free Documentation License	111
1.	APPLICABILITY AND DEFINITIONS	111
2.	VERBATIM COPYING	112
3.	COPYING IN QUANTITY	113
4.	MODIFICATIONS	113
5.	COMBINING DOCUMENTS	114
6.	COLLECTIONS OF DOCUMENTS	115
7.	AGGREGATION WITH INDEPENDENT WORKS	115
8.	TRANSLATION	115
9.	TERMINATION	115
10.	FUTURE REVISIONS OF THIS LICENSE	116
	ADDENDUM: How to use this License for your documents	116
	Glossary of Terms and Acronyms	117
	Index	119

Preface

The audience for this document is software developers that have to use directly the services of XtratuM. The reader is expected to have strong knowledge of the LEON3 (SPARC v8) architecture and experience in programming device drivers. It is also advisable to have some knowledge of the ARINC-653 and related standards.

Typographical conventions

The following font conventions are used in this document:

- **typewriter:** used in assembly and C code examples, and to show the output of commands.
- *italic:* used to introduce new terms.
- **bold face:** used to emphasize or highlight a word or paragraph.

Code

Code examples are printed inside a box like this:

```
static inline void XM_sparcv8_set_psr(xm_u32_t flags) {
    __asm__ __volatile__ ("mov "TO_STR(sparcv8_set_psr_nr)", %%o0\n\t" \
        "mov %0, %%o1\n\t" \
        __DO_XMAHC :: "r"(flags) : "o0", "o1");
}
```

Listing 1: Sample code

Caution sign

The caution sign stresses information that is critical to the integrity or continuity of the system.



Support

Alfons Crespo
Universidad Politecnica de Valencia
Instituto de Automatica e Informtica Industrial
Camino de vera s/n

CP: 46022
Valencia, Spain

The official XtratuM web site is: <http://www.xtratum.org>

Acknowledgements

Porting to LEON3 with MMU has been done in the frame of the “ESA Project AO5829 Securely Partitioning Spacecraft Computing Resources” led by Astrium EADS.

Chapter 1

Introduction

This document describes the XtratuM hypervisor, and how to write applications to be executed as XtratuM partitions.

A hypervisor is a layer of software that provides one or more virtual execution environments for partitions. Although virtualisation concepts has been employed since the 60's (IBM 360), the application of these concepts to the server, desktop, and recently the embedded and real-time computer segments, is a relatively new. There have been some attempts, in the desktop and server markets, to standardise “how” an hypervisor should operate, but the research and the market is not mature enough. In fact, there is still not a common agreement on the terms used to refer to some of the new objects introduced. Check the glossary A.2 for the exact meaning of the terms used in this document.

In the case of embedded systems and, in particular, in avionics, the ARINC-653 standard defines a partitioning system. Although the ARINC-653 standard was not designed to describe how a hypervisor has to operate, some parts of the APEX model of ARINC-653 are quite close to the functionality provided by a hypervisor.

During the porting of XtratuM to the LEON2 and LEON3 processors, we have also adapted the XtratuM API and internal operations to resemble ARINC-653 standard. It is not our intention to convert XtratuM in an ARINC-653 compliant system. ARINC-653 relies on the idea of a “*separation kernel*”, which basically consists in extending and enforcing the isolation between processes or a group of processes. ARINC-653 defines both the API and operation of the partitions, but also how the threads or processes are managed inside each partition. It provides an complete APEX.

In a bare hypervisor, and in particular in XtratuM, a partition is a *virtual computer* rather than a group of strongly isolated processes. When multi-threading (or tasking) support is needed in a partition, then an operating system or a run-time support library has to provide support to the application threads. In fact, it is possible to run a different operating system on each XtratuM partition.

It is important to point out that XtratuM is a bare-metal hypervisor with extended capabilities for highly critical systems. XtratuM provides a raw (close to the native hardware) virtual execution environment, rather than a full featured one. Therefore, **although XtratuM by itself can not be compatible with the ARINC-653 standard, the philosophy of the ARINC-653 has been employed when applicable.**

This document is organised as follows: chapter 2 describes the XtratuM architecture describing how the partitions are organised and scheduled; also, an overview of the XtratuM services is presented.

Chapter 3 outlines the development process on XtratuM: roles, elements, etc.

Chapter 4 describes the compilation process, which involves several steps to finally obtain a binary code which has to be loaded in the embedded system.



The goal of chapter 5 is to provide a view of the API provided by XtratuM to develop applications to be executed as partitions. The chapter puts more emphasis in the development of bare-applications than applications running on a real-time operating system.

Chapter 6 deals with the concrete structure and internal formats of the different components involved in the system development: system image, partition format, partition tables. The chapter ends with the description of the hypercall mechanism.

Chapter 7 and 8 detail the booting process and the configuration elements of the system, respectively. Finally, chapter 8 provides information of the preliminar tools developed to analyse system configuration schemas (XML format) and generate the appropriate internal structures to configure XtratuM for a specific payload.

1.1 History

The term XtratuM derives from the word “stratum”. In geology and related fields it means:

Layer of rock or soil with internally consistent characteristics that distinguishes it from contiguous layers.

In order to stress the tight relation with Linux and the open source the “S” was replaced by “X”. XtratuM would be the first layer of software (the one closer to the hardware), which provides a rock solid basis for the rest of the system.

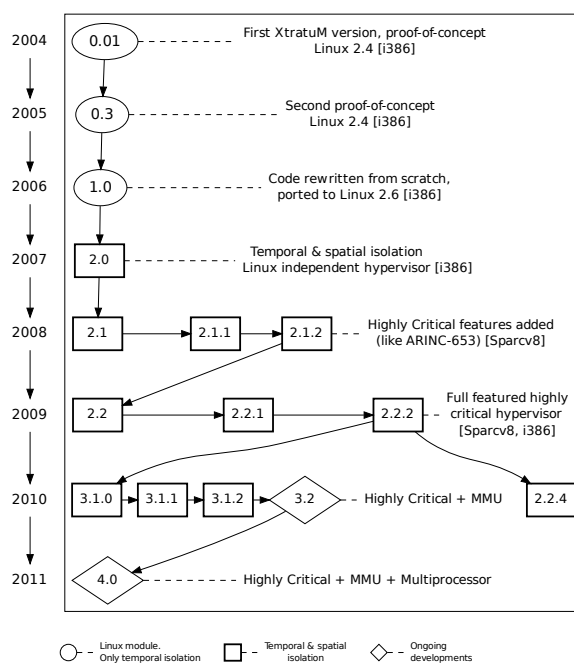


Figure 1.1: XtratuM evolution.

The first version of XtratuM (1.0) was initially developed to meet the requirements of a hard real-time system. The main goal of XtratuM 1.0 was to guarantee the temporal constraints for the real-time partitions. Other characteristics of this version are:

- The first partition shall be a modified version of Linux.
- Partition code has to be loaded dynamically.

- There is not a strong memory isolation between partitions. 55
- Linux is executed in processor supervisor mode.
- Linux is responsible of booting the computer.
- Fixed priority partition scheduling.

XtratuM 2.0 was a completely new redesign and implementation. This new version had nothing in common with the first one but the name. It was a truly hypervisor with both, spatial and temporal isolation. This version was developed for the x86 architecture but never released. 60

XtratuM 2.1 was the first porting to the LEON2 processor, and several safety critical features were added. Just to mention the most relevant features:

- Bare metal hypervisor.
- Employs para-virtualisation techniques. 65
- A hypervisor designed for embedded systems: some devices can be directly managed by a designated partition.
- Strong temporal isolation: fixed cyclic scheduler.
- Strong spatial isolation: all partitions are executed in processor user mode, and do not share memory. 70
- Resource allocation via a configuration table.
- Robust communication mechanisms (ARINC sampling and queuing ports).

Version 2.1 was a prototype to evaluate the capabilities of the LEON2 processor to support a hypervisor system.

XtratuM 2.2 was a more mature hypervisor on the LEON2 processor. This version has most of the final functionality. 75

The current development version is 3.1, which contains MMU support. Version 3.1 is still under active development. The first stable version with MMU support will be named 3.2.

In what follows, the name XtratuM will be used to refer to the version 3.1 and eventually to 3.2 of XtratuM. 80

This page is intentionally left blank.

Chapter 2

XtratuM Architecture

This chapter introduces the architecture of XtratuM.

The concept of partitioned software architectures was developed to address security and safety issues. The central design criteria involves isolating modules of the system into *partitions*. Temporal and spatial isolation are the key aspects in a partitioned system. Based on this approach, the Integrated Modular Avionics (IMA) is a solution that allowed the Aeronautic Industry to manage the increment of the functionalities of the software maintaining the level of efficiency.

85

XtratuM is a bare-metal hypervisor that has been designed to achieve temporal and spatial partitioning for safety critical applications. Figure 2.1 shows the complete architecture.

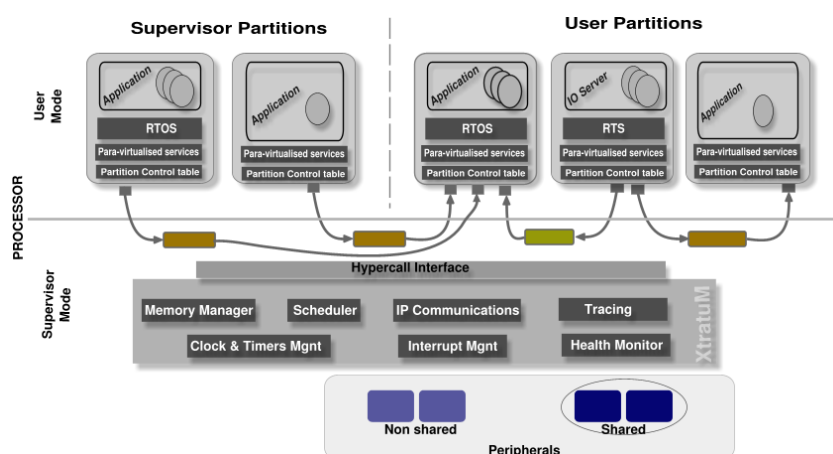


Figure 2.1: XtratuM architecture.

The main components of this architecture are:

- Hypervisor: XtratuM provides virtualisation services to partitions. It is executed in supervisor processor mode and virtualises the CPU, memory, interrupts, and some specific peripherals. The internal XtratuM architecture includes the following components:
 - Memory management: XtratuM provides a memory model for the partitions enforcing the spatial isolation. It uses the hardware mechanisms to guarantee the isolation.
 - Scheduling: Partitions are scheduled using a cyclic scheduling policy.
 - Interrupt management: Interrupts are handled by XtratuM and, depending on the interrupt nature, propagated to the partitions. XtratuM provides an interrupt model to the partitions that extends the concept of processor interrupts by adding 32 additional interrupt numbers.

90

95

– Clock and timer management:

– IP communication: Inter-partition communication is related with the communications between two partitions or between a partition and the hypervisor. XtratuM implements a message passing model which highly resembles the one defined in the ARINC-653. A communication channel is the logical path between one source and one or more destinations. Two basic transfer modes are provided: sampling and queuing. Partitions can access to channels through access points named ports. The hypervisor is responsible for encapsulating and transporting messages.

– Health monitor: The health monitor is the part of XtratuM that detects and reacts to anomalous events or states. The purpose of the HM is to discover the errors at an early stage and try to solve or confine the faulting subsystem in order to avoid or reduce the possible consequences.

– Tracing facilities: XtratuM provides a mechanism to store and retrieve the traces generated by partitions and XtratuM itself. Traces can be used for debugging, during the development phase of the application, but also to log relevant events or states during the production phase.

- API: Defines the para-virtualised services provided by XtratuM. The access to these services is provided through *hypercalls*.

- Partitions: A partition is an execution environment managed by the hypervisor which uses the virtualised services. Each partition consists of one or more concurrent processes (implemented by the operating system of each partition), that share access to processor resources based upon the requirements of the application. The partition code can be: an application compiled to be executed on a bare-machine; a real-time operating system (or runtime support) and its applications; or a general purpose operating system and its applications.

Partitions need to be *virtualised* to be executed on top of a hypervisor. Depending on the type of execution environment, the virtualisation implications in each case can be summarised as:

Bare application : The application has to be virtualised using the services provided by XtratuM. The application is designed to run directly on the hardware and the hardware must be aware of this fact.

Operating system application : When the application runs on top of a (real-time) operating system, it uses the services provided by the operating system and does not need to be virtualised. However, the operating system has to deal with the virtualisation and be virtualised (ported on top of XtratuM).

2.1 System operation

The system states and its transitions are shown in figure 2.2.

At boot time, the resident software loads the image of XtratuM in main memory and transfers the control to the entry point of XtratuM. The period of time between starting from the entry point, to the execution of the first partition is defined as **boot** state. In this state, the scheduler is not enabled and the partitions are not executed (see chapter 7).

At the end of the boot sequence, the hypervisor is ready to start executing partition code. The system changes to **normal** state and the scheduling plan is started. Changing from boot to normal state is performed automatically (the last action of the set up procedure).

The system can switch to **halt** state by the health monitoring system in response to a detected error or by a *system partition* invoking the service `XM.halt_system()`. In the halt state: the scheduler is

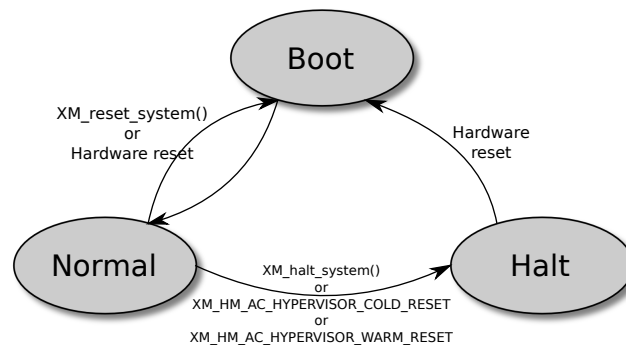


Figure 2.2: System states and transitions.

disabled, the hardware interrupts are disabled, and the processor enters in an endless loop. The only way to exit from this state is via an external hardware reset.

It is possible to perform a warm or cold (hardware reset) system reset by using the hypercall (see `XM_reset_system()`). On a warm reset, the system increments the reset counter, and a reset value is passed to the new rebooted system. On a cold reset, no information about the state of the system is passed to the new rebooted system.

145

2.2 Partition operation

Once XtratuM is in normal state, partitions are started. The partition's states and transitions are shown in figure 2.3.

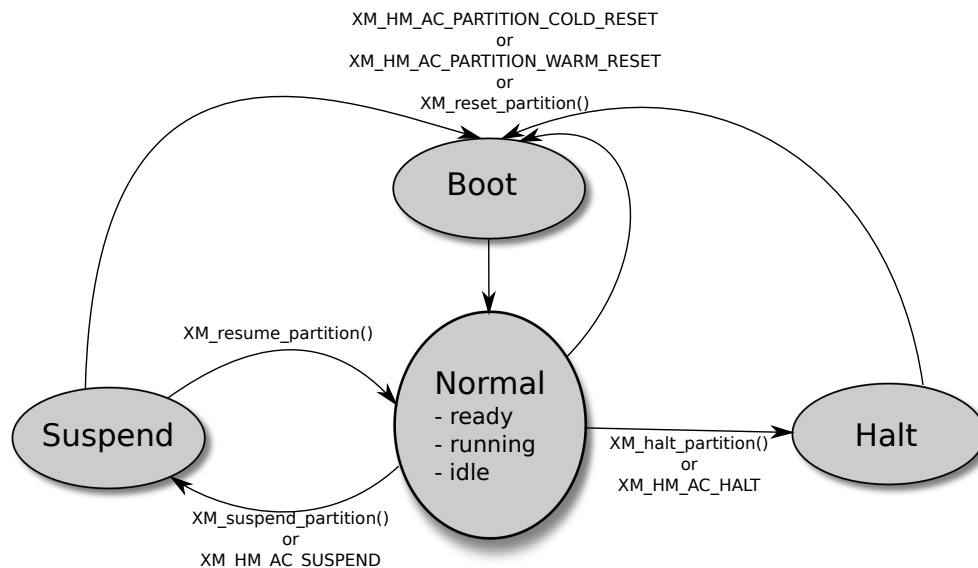


Figure 2.3: Partition states and transitions.

On start-up each partition is in boot state. It has to prepare the virtual machine to be able to run the

150

applications¹: it sets up a standard execution environment (that is, initialises a correct stack and sets up the virtual processor control registers), creates the communication ports, requests the hardware devices (I/O ports and interrupt lines), etc., that it will use. Once the operating system has been initialised, the partition changes to normal mode.

155 The partition receives information from XtratuM about the previous executions, if any.

From the point of view of the hypervisor, there is no difference between the boot state and the normal state. In both states the partition is scheduled according to the fixed plan, and has the same capabilities. Although not mandatory, it is recommended that the partition emits a partition's state-change event when changing from boot to normal state.

160 The normal state is subdivided in three sub-states:

Ready The partition is ready to execute code, but is not scheduled because it is not in its time slot.

Running The partition is being executed by the processor.

Idle If the partition does not want to use the processor during its allocated time slot, it can relinquish the processor, and waits for an interrupt or for the next time slot (see `XM_idle_self()`).

165 A partition can be moved to the halt state by itself or by a system partition. In the halt state, the partition is not executed by the scheduler and the time slot allocate to it is left idle (it is not allocated to other partitions). All the resources allocated to the partition are released. It is not possible to return to normal state.

170 In suspend state, a partition will not be scheduled and interrupts are not delivered. Interrupts raised while in suspended state are left pending. If the partition returns to normal state, then pending interrupts are delivered to the partition. The partition can return to the normal state if requested by a system partition by calling `XM_resume_partition()` hypercall.

2.3 System partitions

XtratuM defines two types of partitions: *normal* and *system*. System partitions are allowed to manage and monitor the state of the system and other partitions. Some hypercalls cannot be called by a normal partition or have restricted functionality.

175 Note that system rights are related to the capability to manage the system, and not to the capability to access directly to the native hardware or to break the isolation: a system partition is scheduled as a normal partition; and it can only use the resources allocated to it in the configuration file.

180 Table 2.1 shows the list of hypercalls reserved for system partitions. A hypercall labeled as “partial” indicates that a normal partition can invoke it if a system reserved service is not requested.

A partition has system capabilities if the `/System_Description/Partition_Table/Partition/@flags` attribute contains the flag “system” in the XML configuration file.

2.4 Names and identifiers

185 Each partition is globally identified by an unique number *id*. Partition identifiers are assigned by the integrator in the `XM_CF` file. XtratuM uses this number to refer to partitions. System partitions use partition identifiers to refer to the target partition. There “C” macro `XM_PARTITION_SELF` can be used by a partition to refer to itself.

¹We will consider that the partition code is composed of an operating system and a set of applications.

Hypercall	System
<code>XM_get_partition_status</code>	Yes
<code>XM_get_plan_status</code>	Yes
<code>XM_get_system_status</code>	Yes
<code>XM_halt_partition</code>	Partial
<code>XM_halt_system</code>	Yes
<code>XM_hm_open</code>	Yes
<code>XM_hm_read</code>	Yes
<code>XM_hm_seek</code>	Yes
<code>XM_hm_status</code>	Yes
<code>XM_memory_copy</code>	Partial
<code>XM_reset_partition</code>	Partial
<code>XM_reset_system</code>	Yes
<code>XM_resume_partition</code>	Yes
<code>XM_set_plan</code>	Yes
<code>XM_shutdown_partition</code>	Partial
<code>XM_suspend_partition</code>	Partial
<code>XM_trace_open</code>	Yes
<code>XM_trace_read</code>	Yes
<code>XM_trace_seek</code>	Yes
<code>XM_trace_status</code>	Yes

Table 2.1: List of system reserved hypercalls.

These id's are used internally as indexes to the corresponding data structures². The first "id" of each object group shall start in zero and the next id's shall be consecutive. It is mandatory to follow this ordering in the XM_CF file.

The attribute *name* of a partition is a human readable string. This string shall contain only the following set of characters: upper and lower case letters, numbers and the underscore symbol. It is advisable not to use the same name on different partitions. A system partition can get the name of another partition by consulting the status object of the target partition.

In order to avoid name collisions, all the public symbols of XtratuM contain the prefix "xm". Therefore, the prefix "xm", both in upper and lower case, is reserved.

2.5 Partition scheduling

XtratuM schedules partitions in a fixed, cyclic basis (ARINC-653 scheduling policy). This policy ensures that one partition cannot utilise the processor for longer than intended to the detriment of the other partitions. The set of *time slots* allocated to each partition are defined in the XM_CF configuration during the design phase. Each partition is scheduled for a time slot defined as a starting time and a duration. Within a time slot, XtratuM allocates the processor to the partition.

If there are several concurrent activities in the partition, the partition shall implement its own scheduling algorithm. This two-level scheduling scheme is known as *hierarchical scheduling*. XtratuM is not aware of the scheduling policy used internally on each partition.

In general, a cyclic plan consists in a major time frame (MAF) which is periodically repeated. The MAF is typically defined as the least common multiple of the periods of the partitions (or the periods of the threads of each partition, if any).

²For efficiency and simplicity reasons.

	Name	Period	WCET	Util %
Partition 1	System Mngmt	100	20	20
Partition 2	Flight Control	100	10	10
Partition 3	Flight Mngmt	100	30	30
Partition 4	IO Processing	100	20	20
Partition 5	IHVM	200	20	10

(a) Partition set.

	Start	Dur.	Start	Dur.	Start	Dur.	Start	Dur.
Partition 1	0	20	100	20				
Partition 2	20	10	120	10				
Partition 3	40	30	140	30				
Partition 4	30	10	70	10	130	10	170	10
Partition 5	180	20						

(b) Detailed execution plan.

Table 2.2: Partition definition.

For instance, consider the partition set of figure 2.2a, its hyper-period is 200 time units (milliseconds) and has a CPU utilisation of the 90%. The execution chronogram is depicted in figure 2.4. One of the possible cyclic scheduling plan can be described, in terms of start time and duration, as it is shown in the table 2.2b.

210

This plan has to be specified in the configuration file. An XML file describing this schedule is shown below.

```

<CyclicPlanTable>
  <Plan majorFrame="1s">
    <Slot duration="20ms" id="0" partitionId="0" start="0ms"/>
    <Slot duration="10ms" id="1" partitionId="1" start="20ms"/>
    <Slot duration="10ms" id="2" partitionId="0" start="30ms"/>
    <Slot duration="30ms" id="3" partitionId="2" start="40ms"/>
    <Slot duration="10ms" id="4" partitionId="1" start="70ms"/>
    <Slot duration="20ms" id="5" partitionId="0" start="100ms"/>
    <Slot duration="10ms" id="6" partitionId="1" start="120ms"/>
    <Slot duration="10ms" id="7" partitionId="0" start="130ms"/>
    <Slot duration="30ms" id="8" partitionId="2" start="140ms"/>
    <Slot duration="10ms" id="9" partitionId="1" start="170ms"/>
    <Slot duration="20ms" id="10" partitionId="0" start="180ms"/>
  </Plan>
</CyclicPlanTable>
    
```

Listing 2.1: /user/examples/sched_events/xm.cf.sparcv8.xml

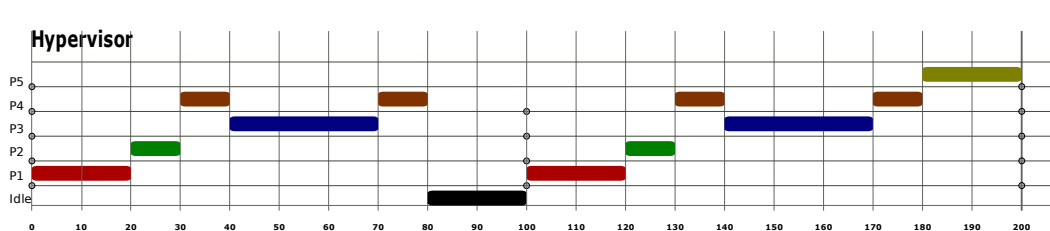


Figure 2.4: Scheduling example.

One important aspect in the design of the XtratuM hypervisor scheduler is the consideration of the

overhead caused by the partition's context switch. Figure 2.5 shows the implications of this issue. Subfigure 2.5a shows the context switch between partitions 1 and 2. To execute the partition, XtratuM saves the partition 1's context and loads the partition 2's context. 215

XtratuM scheduling design tries to adjust as much as possible the beginning of the execution to the specified starting time of the slot. To do that, when a slot is scheduled, XtratuM programs a timer with the duration of the slot minus the temporal cost of the complete context switch (load and save the context). Subfigure 2.5b shows this situation. 220

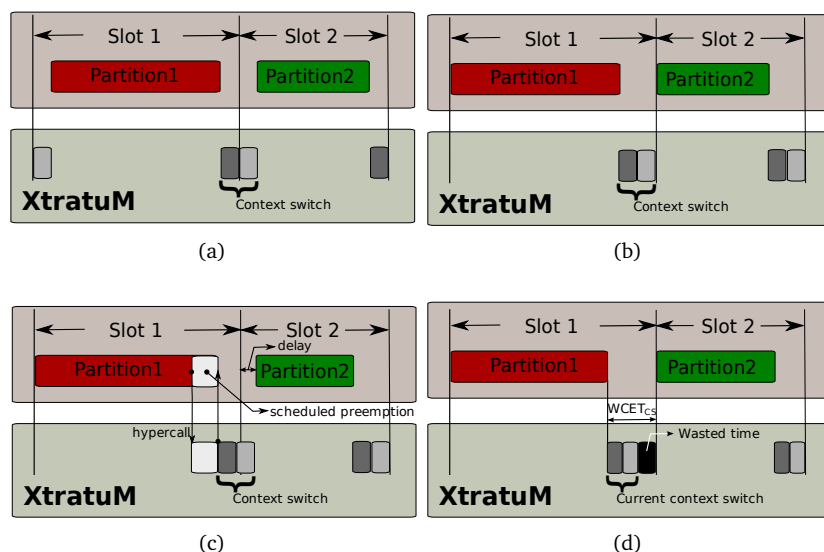


Figure 2.5: XtratuM context switch analysis.

However, the scenario depicted in subfigure 2.5c can occur. In this case, just before the duration timer expiration the partition invokes a hypercall. When the hypercall finishes, the timer interrupt is detected and the context switch is done at that time. This situation can introduce some small delay in the beginning of partition of the next scheduling time slot.

Figure 2.5d details what should be the value of the considered cost of the context switch. If the duration of the context switch is assumed as the worst case execution time of the context switch ($WCET_{CS}$), a situation like the one shown in figure 2.5d may happen. In this example, the cost of the context switch is less than its $WCET_{CS}$ and, as consequence, an idle time has to be introduced to start the execution of the partition at the specified time. 225

XtratuM copes this situation by implementing the following algorithm: 230

- When a partition is scheduled, a timer (Scheduler Timer, ST) is armed with a value that considers the absolute start time of the next time slot, and the best case execution time of the context switch ($BCET_{CS}$).
- Two situations can introduce a small delay to the effective starting of the slot:
 1. The actual cost of the context switch is larger than the $BCET_{CS}$. In this case, the execution will start with a delay that is $WCET_{CS} - BCET_{CS}$. 235
 2. The ST expires while a hypercall is under execution. XtratuM will carry out the context switch when the current hypercall is finished, which delays the context switch. The worst case situation corresponds to the hypercall with longer execution time: $WCET_{HP}$.

Both previous situations can occur simultaneously. So, the worst case delay can be estimated as $(WCET_{CS} - BCET_{CS}) + WCET_{HP}$. 240

The cost of the context switch (both: $WCET_{CS}$ and $BCET_{CS}$) and all hypercalls have been evaluated and identified the worst case situation. In the document “*Volume 3: Testing and Evaluation*” it is provided a deep analysis of the hypercalls. The integrator must consider the worst case execution time of the used hypercalls and the partition context switch to forecast the slot duration considering the hypercalls used in the partition and the XtratuM configuration parameters.

2.5.1 Multiple scheduling plans

In some cases, a single scheduling plan may be too restrictive. For example:

- Depending on the guest operating system, the initialisation can require a certain amount of time and can vary significantly. If there is a single plan, the initialisation of each partition can require different number of slots due to the fact that the slot duration has been designed considering the operational mode. This implies that a partition can be executing operational work whereas other are still initialising its data.
- The system can require to execute some maintenance operations. These operation can require other resource allocation than the operational mode.

In order to deal with these issues, XtratuM provides multiple scheduling plans that allows to reallocate the timing resources (the processor) in a controlled way. In the scheduling theory this process is known as mode changes. Figure 2.6 shows how the modes have been considered in the XtratuM scheduling.

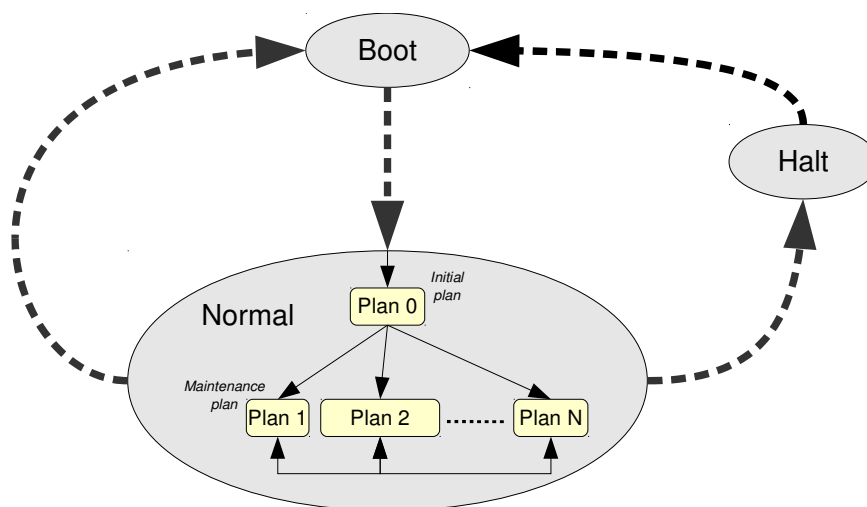


Figure 2.6: Scheduling modes.

The scheduler (and so the plans) are only active while the system is in *normal* mode. Plans are defined in the XM_CF file and identified with a number. Some plans are reserved or have a special meaning:

Plan 0: *Initial plan.* The system executes this plan after a system reset. The system will be in plan 0 until a plan change is requested.

It is not legal to switch back to this plan. That is, this plan is only executed as a consequence of a system reset (software or hardware).

Plan 1: *Maintenance plan.* This plan can be activated in two ways:

- As the a result of the health monitoring action `XM_HM_AC_SWITCH_TO_MAINTENANCE`. The plan switch is done immediately.

- Requested from a system partition. The plan switch occurs at the end the current plan.

It is advisable to allocate the first slot of this plan to a system partition, in order to start the maintenance activity as soon as possible after the plan switch. Once the maintenance activities has been completed, it is responsibility of a system partition to switch to another plan (if needed).

270

A system partition can also request a switch to this.

Plan x (x>1): Any plan greater than 1 is used defined. A system partition can switch to any defined plan at any time.

Switching scheduling plans

When a plan switch is requested by a system partition (through a hypercall), the plan switch is not synchronous; all the slots of the current plan will be completed, and the new plan will be started at the end of the current one.

275

The plan switch that occurs as a consequence of the `XM_HM_AC_SWITCH_TO_MAINTENANCE` action is synchronous. The current slot is terminated, and the Plan 1 is started immediately.

2.6 Inter-partition communications (IPC)

Inter-partition communications are related with the communications between two partitions. XtratuM implements a message passing model which highly resembles the one defined in the ARINC-653 standard. A message is a variable³ block of data. A message is sent from a partition source to one or more partitions' destinations. The data of a message is transparent to the message passing system.

280

A communication channel is the logical path between one source and one or more destinations. Partitions can access to channels through access points named ports. The hypervisor is responsible for encapsulating and transporting messages that has to arrive to the destination(s) unchanged. At the partition level, messages are atomic entities i.e., either the whole message is received or nothing is received. Partition developers are responsible for agreeing on the format (data types, endianness, padding, etc.).

285

Channels, ports, maximum message sizes and maximum number of messages (queuing ports) are entirely defined in the configuration files (see section 8).

290

XtratuM provides two basic transfer modes: *sampling* and *queuing*.

Sampling port: It provides support for broadcast, multicast and unicast messages. No queuing is supported in this mode. A message remains in the source port until it is transmitted through the channel or it is overwritten by a new occurrence of the message, whatever occurs first. Each new instance of a message overwrites the current message when it reaches a destination port, and remains there until it is overwritten. This allows the destination partitions to access the latest message.

295

A partition's write operation on a specified port is supported by `XM_write_sampling_message()` hypercall. This hypercall copies the message into an internal XtratuM buffer. Partitions can read the message by using `XM_read_sampling_message()` which returns the last message written in the buffer. XtratuM copies the message to the partition space.

300

Any operation on a sampling port is non-blocking: a source partition can always write into the buffer and the destination partition/s can read the last written message.

The channel has an optional configuration attribute named `@refreshPeriod`. This attribute defines the maximum time that the data written in the channel is considered "valid". Messages older

305

³XtratuM defines the maximum length of a message.

than the valid period are marked as invalid. When a message is read, a bit is set accordingly to the valid state of the message.

Queueing port: It provides support for buffered unicast communication between partitions. Each port has associated a queue where messages are buffered until they are delivered to the destination partition. Messages are delivered in FIFO order.

Sending and receiving messages are performed by two hypercalls: `XM_send_queueing_message()` and `XM_receive_queueing_message()`, respectively. XtratuM implements a classical producer-consumer circular buffer without blocking. The sending operation writes the message from partition space into the circular buffer and the receive one performs a copy from the XtratuM circular buffer into the destination memory.



If the requested operation cannot be completed because the buffer is full (when trying to send a message) or empty (when attempting to receive a message), then **the operation returns immediately with the corresponding error**. The partition's code is responsible for retrying the operation later.

In order to optimise partition's resources and reduce the performance loss caused by polling the state of the port. XtratuM triggers an extended interrupt when the a new message is written/sent to a port. Since there is only one single interrupt line to notify for incoming messages, on the reception of the interrupt, the partition code has to determine which port or ports are ready to perform the operation. XtratuM maintains a bitmap in the Partition Control Table to inform about the state of each port. A "1" in the corresponding entry indicates that the requested operation can be performed.

When a new message is available in the channel, XtratuM triggers an extended interrupt to the destination(s).

2.7 Health monitor (HM)

The health monitor is the part of XtratuM that detects and reacts to anomalous events or states. The purpose of the HM is to discover the errors at an early stage and try to solve or confine the faulting subsystem in order to avoid a failure or reduce the possible consequences.

It is important to clearly understand the difference between 1) an incorrect operation (instruction, function, application, peripheral, etc.) which is handled by the normal control flow of the software, and 2) an incorrect behaviour which affects the normal flow of control in a way not considered by the developer or which can not be handled in the current scope.

An example of the first kind of errors is when the `malloc()` function returns a null pointer when there are not enough memory to attend the request. This error is typically handled by the program by checking the return value. An attempt to execute a undefined instruction (processor instruction) may not be properly handled by the program that attempted to execute it.

The XtratuM health monitoring system will manage those faults that cannot, or should not, be managed at the *scope* where the fault occurs.

The XtratuM HM system is composed of four logical blocks:

HM event detection:

to detect abnormal states, using logical probes in the XtratuM code.

HM actions:

a set of predefined actions to recover the fault or confine the error.

HM configuration:

to bind the occurrence of each HM event with the appropriate HM action.

HM notification:

to report the occurrence of the HM events.

Since HM events are, by definition, the result of a non-expected behaviour of the system, it may be difficult to clearly determine which is the original cause of the fault, and so, which is the best way to handle the problem. XtratuM provides a set of “coarse grain” actions (see section 2.7.2) that can be employed at the first stage, right when the fault is detected. Although XtratuM implements a default action for each HM event, the integrator can map an HM action to each HM event using the XML configuration file.

Once the defined HM action is carried out by XtratuM, a HM notification message is stored in the HM log stream (if the hm event is marked to generate a log). A system partition can then read those log messages and perform a more advanced error handling. As an example of what can be implemented:

1. Configure the hm action to stop the faulting partition, and log the event.
2. The system partition can resume an alternate one, a redundant dormant partition, which can be implemented by another developer team to achieve diversity.

Since the differences between *fault*⁴ and *error*⁵ are so subtle and subjective, we will use both terms to refer to the original reason of an incorrect state.

The XtratuM health monitoring subsystem defines four different execution scopes, depending on which part of the system has been initially affected:

1. Process scope: Partition process or thread.
2. Partition scope: Partition operating system or run-time support.
3. Hypervisor scope: XtratuM code.
4. Board scope: Resident software (BIOS, BOOT ROM or firmware).

The scope⁶ where an HM event should be managed has to be greater than the scope where it was “believed” to be produced.

There is not a clear and unique scope for each HM event. Therefore the same HM event may be handled at different scopes. For example, fetching an illegal instruction is considered hypervisor scope if it happens when while XtratuM is executing; and partition level if the event is raised while a partition is running.

XtratuM tries to determine the most likely scope target, and the delivers the HM to the corresponding upper scope.

Note that although in the LEON2 version of XtratuM there is no distinction between the first and second scopes, it is important to consider that there are two different parts in the partition’s code: user applications, and operating system. Therefore, it is consistent to deliver the first scope of HM events, caused by a process or thread, to the second scope.

2.7.1 HM Events

There are three sources of HM events:

⁴Fault: What is believed to be the original reason that caused an error.

⁵Error: The manifestation of a fault.

⁶The term **level** is used in the ARINC-653 standard to refer to this idea

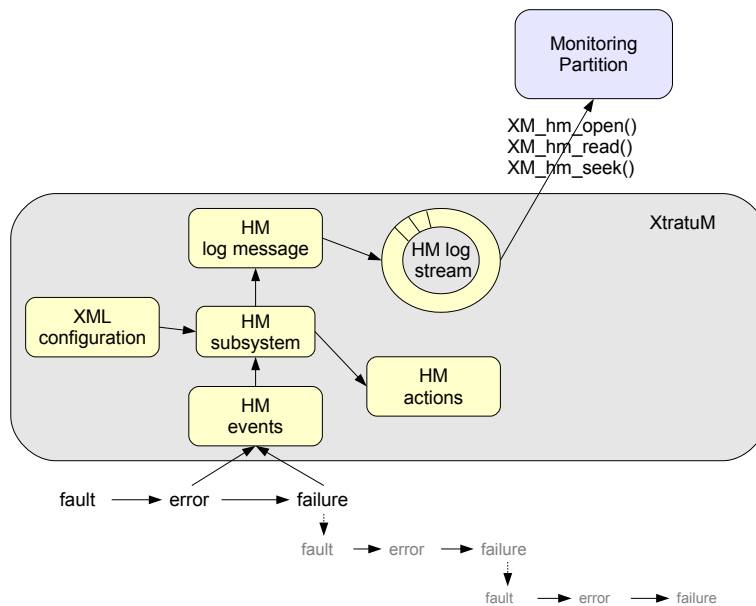


Figure 2.7: Health monitoring overview.

- Events caused by abnormal hardware behaviour. These events are notified to XtratuM via processor traps. Most of the processor exceptions are managed as health monitoring events.
- Events detected and triggered by partition code. These events are usually related to checks or assertions on the code of the partitions. **Health monitoring events raised by partitions are a special type of tracing message** (see sections 2.10). Highly critical tracing messages are considered as HM events.
- Events triggered by XtratuM. Caused by a violation of a sanity check performed by XtratuM on its internal state or the state of a partition.

When the HM event is detected, the relevant information (error scope, offending partition id, memory address, faulting device, etc.) is gathered and used to select the appropriate HM action.

2.7.2 HM Actions

Once an HM event is raised, XtratuM has to react quickly to the event. The set of configurable HM actions are listed in the next table:

Action	Description
XM_HM_AC_IGNORE	No action is performed.
XM_HM_AC_SHUTDOWN	The shutdown extended interrupt is sent to the failing partition.
XM_HM_AC_COLD_RESET	The failing partition/processor is cold reset.
XM_HM_AC_WARM_RESET	The failing partition/processor is warm reset.
XM_HM_AC_SUSPEND	The failing partition is suspended.
XM_HM_AC_HALT	The failing partition/processor is halted.
XM_HM_AC_PROPAGATE	No action is performed by XtratuM. The event is redirected to the partition as a virtual trap.
XM_HM_AC_SWITCH_TO_MAINTENANCE	The current scheduling plan is switched to the maintenance one.

2.7.3 HM Configuration

There are two tables to bind the HM events with the desired handling actions:

XtratuM HM table: which defines the actions for those events that has to be managed at system or hypervisor scope;

Partition HM table: which defined the actions for those events that has to be managed at hypervisor or partition scope;

400

Note that the same HM event can be binded with different recovery actions in each partition HM table and in the XtratuM HM table.

The HM system can be configured to send an HM message after the execution of the HM action. It is possible to select whether an HM event is logger or not. See the chapter 8.

2.7.4 HM notification

The log events generated by the hm system (those event that are configured to generate a log) are stored in the device configured in the XM.CF configuration file.

405

In the case that the logs are stored in a log stream, then they can be retrieved by system partitions using the XM.hm.X services.

The maximum number of messages on is configured in the XtratuM source code (see the section 8.1).

Health monitoring log messages are fixed length messages defined as follows:

```
typedef struct {
    xm_u32_t eventId:13, system:1, reserved:1, validCpuCtxt:1, moduleId:8,
    partitionId:8;
    struct cpuCtxt cpuCtxt;
#define HM_PAYLOAD_WORDS 5
    xm_u32_t word[HM_PAYLOAD_WORDS]; /* Payload */
    xmTime_t timeStamp;
} xmHmLog_t;
```

Listing 2.2: /core/include/objects/hm.h

eventId: Identifies the event that caused this log.

410

system: Set if the error was raised while executing XtratuM code.

moduleId: In the case of events raised by a partition (as a consequence of a high critical trace message), this field is a copy of the field with the same name of the trace message.

partitionId: The Id attribute of the partition that may caused the event.

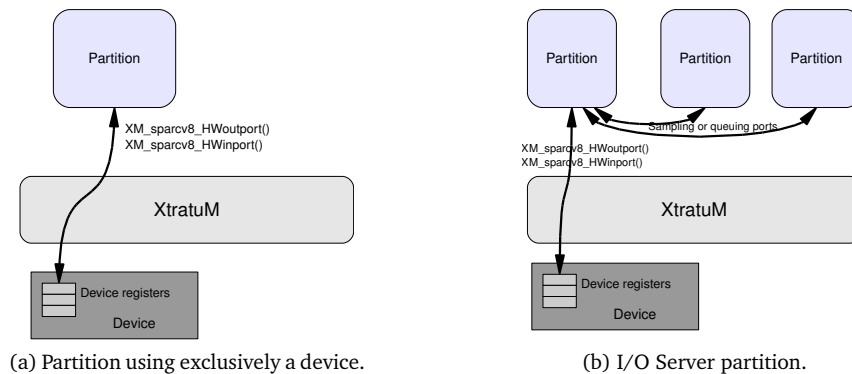
word: Event specific information.

415

timeStamp: A time stamp of when the event was detected.

2.8 Access to devices

A partition, using exclusively a device (peripheral), can access the device through the device driver implemented in the partition (figure 2.8a). The partition is in charge of handling properly the device.



(a) Partition using exclusively a device.

(b) I/O Server partition.

The configuration file has to specify the I/O ports and the interrupt lines that will be used by each partition.

Two partitions cannot use the same the same interrupt line. XtratuM provides a fine grain access control to I/O ports, so that, several partitions can use (read and write) different bits of the the same I/O port. Also, it is possible to define a range of valid values that can be written in a I/O port (see section 5.10).

When a device is used by several partitions, a user implemented I/O server partition (figure 2.8b) may be in charge of the device management. An I/O server partition is a specific partition which accesses and controls the devices attached to it, and exports a set of services via the inter-partitions communication mechanisms provided by XtratuM (sampling or queuing ports), enabling the rest of partitions to make use of the managed peripherals. The policy access (priority, fifo, etc.) is implemented by the I/O server partition.

Note that the I/O server partition is not part of XtratuM. It should, if any, be implemented by the user of XtratuM.

2.9 Traps, interrupts and exceptions

2.9.1 Traps

A **trap** is the mechanism provided by the LEON3 processor to implement the asynchronous transfer of control. When a trap occurs, the processor switches to supervisor mode and unconditionally jumps into a predefined handler.

SPARC v8 defines 256 different trap handlers. The table which contains these handlers is called *trap table*. The address of the trap table is stored in a special processor register (called \$tbr). Both, the \$tbr and the contents of the trap table are exclusively managed by XtratuM. All native traps jump into XtratuM routines.

The trap mechanism is used for several purposes:

Hardware interrupts Used by peripherals to request the attention of the processor.

Software traps Raised by a processor instruction; commonly used to implement the system call mechanism in the operating systems.

Processor exceptions Raised by the processor to inform about a condition that prevent the execution of an instruction. There are basically two kind of exceptions: those caused by the normal operation of the processor (such as register window under/overflow), and those caused by an abnormal situation (such as an memory error).

XtratuM defines 32 new interrupts called *extended interrupts*. These new interrupts are used to inform the partition about XtratuM specific events. Those new traps are vectored at the end of the trap table (entry trap 224 in the case of the SPARC v8). The trap handler raised by a trap can be changed by invoking the `XM_route_irq()` hypercall.

450

Partitions are not allowed to use (read or write) the `$tbr` register. XtratuM implements a *virtual trap table*.

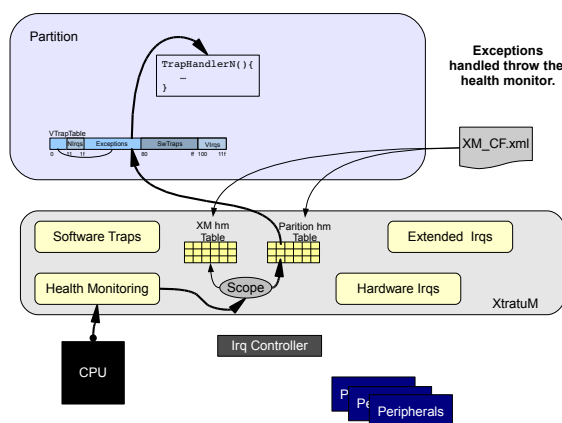


Figure 2.8: Exceptions handled by the health monitoring subsystem.

2.9.2 Interrupts

Although in a fully virtualised environment, a partition should not need to manage hardware interrupt; XtratuM only virtualises those hardware peripherals that may endanger the isolation, but leaves to the partitions to directly manage non-critical devices.

455

In order to properly manage peripherals, a partition needs to:

1. have access to the peripheral control and data registers.
2. be informed about triggered interrupts.
3. be able to block (mask and unmask) the associated interrupt line.

460

A hardware interrupt can only be allocated to one partition (in the `XM_CF` configuration file). The partition can then mask and unmask the hardware line in the native interrupt controller using the `XM_mask_irq()` and `XM_unmask_irq()` functions.

XtratuM extends the concept of processor traps by adding a 32 additional interrupt numbers. This new range is used to inform the partition about events detected or generated by XtratuM.

465

Figure 2.9 shows the sequence from the occurrence of an interrupt to the partition's trap handler.

Partitions shall manage this new set of events in the same way standard traps are. The native trap table of the LEON3 is extended, appending 32 new trap entries, which will be invoked by XtratuM on the occurrence of an event alike a standard LEON3 trap.

2.10 Traces

XtratuM provides a mechanism to store and retrieve the traces generated by partitions and XtratuM itself. Traces can be used for debugging, during the development phase of the application, but also to

470

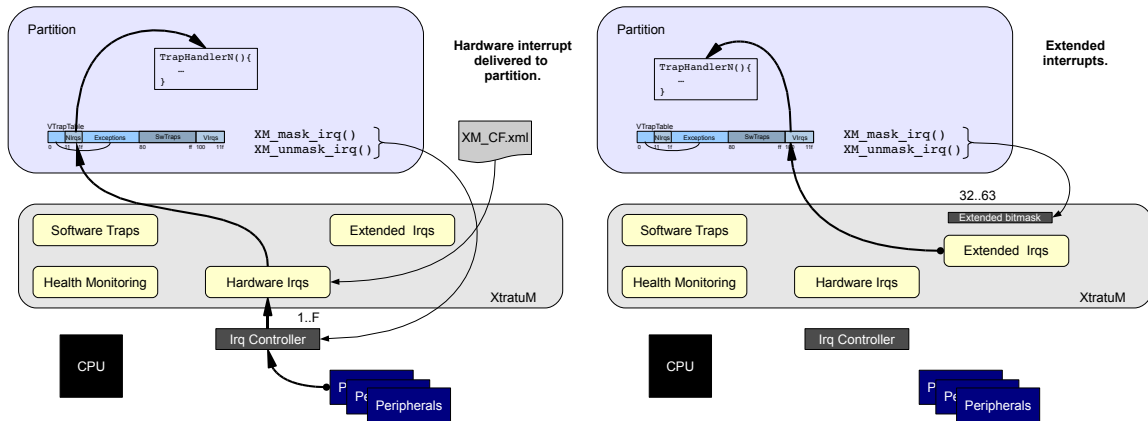


Figure 2.9: Hardware and extended interrupts delivery.

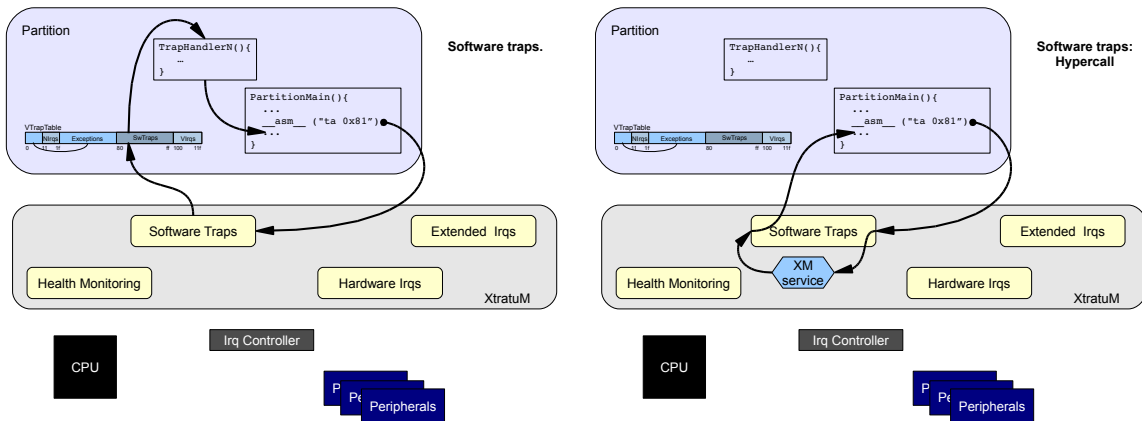


Figure 2.10: Software traps.

log relevant events during the production phase.

In order to enforce resource isolation, each partition (as well as XtratuM) has a dedicated trace log stream to store its own trace messages, which is specified in the @device attribute of the Trace element. Trace is an optional element of XMHypervisor and Partition elements.

475

The hypercall to write a trace message has a parameter (bitmask) used to select the traces messages are stored in the log stream. The integrator can select which trace messages are actually stored in the log stream with the Trace/@bitmask attribute. If the logical and between the value configured in the Partition/Trace/@bitmask and the value of the bitmask parameter of the XM_trace_event() hypercall, then the event is stored, otherwise it is discarded.

480

Figure 2.11 sketches the configuration of the traces. In the example, the traces generated by partition 1 will be stored in the device MemDisk0, which is defined in the Devices section as a memory block device. Only those traces whose least significant bit is set in the bitmask parameter will be recorded.

2.11 Clocks and timers

There are two clocks per partition:

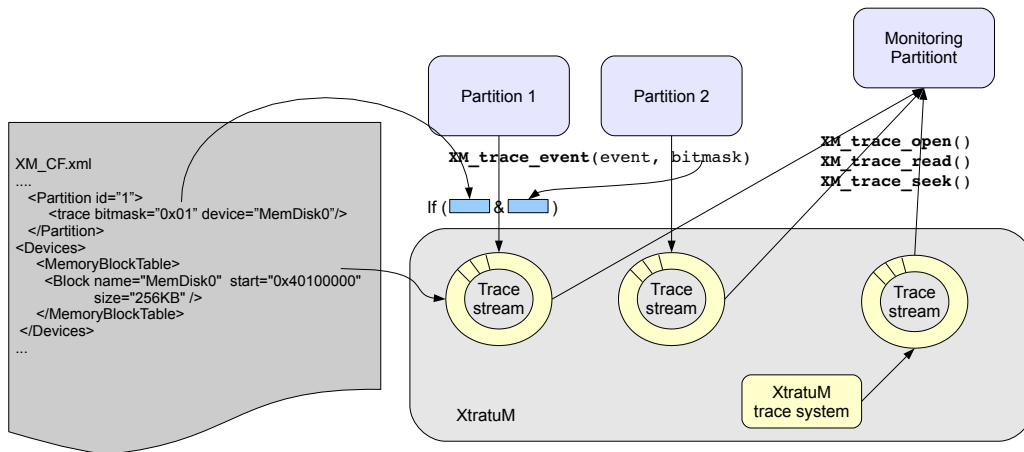


Figure 2.11: Tracing overview.

XM_HW_CLOCK: Associated with the native hardware clock. The resolution is 1μsec.

485

XM_EXEC_CLOCK: Associated with the execution of the partition. This clock only advances while the partition is being executed. It can be used by the partition to detect overruns. This clock relies on the XM_HW_CLOCK and its resolution is also 1μsec.

Only one timer can be armed for each clock.

2.12 Status

Relevant internal information regarding the current state of the XtratuM and the partitions, as well as accounting information is maintained in an internal data structure that can be read by system partitions.

490

This optional feature shall be enabled in the XtratuM source configuration, and then recompile the XtratuM code. **By default it is disabled.** The hypercall is always present; but if not enabled, then XtratuM does not gather statistical information and then some status information fields are undefined. It is enabled in the XtratuM menuconfig: Objects → XM partition status accounting.



495

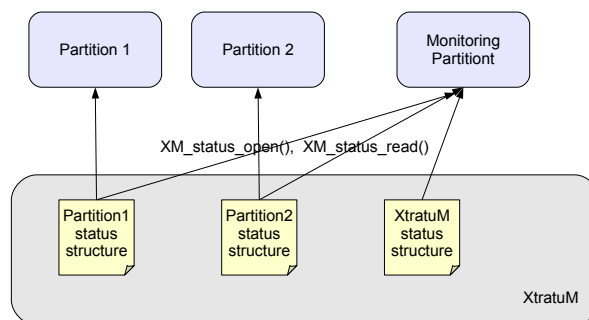


Figure 2.12: Status overview.

2.13 Summary

Next is a brief summary of the ideas and concepts that shall be kept in mind to understand the internal operation of XtratuM and how to use the hypercalls:

- A partition behaves basically as the native computer. Only those services that have been explicitly para-virtualised should be managed in a different way.
- 500 • Partition's code should not be self-modifying.
- Partition's code is always executed with native interrupts enabled.
- Partition's code is not allowed to disable native interrupts, only their own virtual interrupts.
- XtratuM code is non-preemptive. It should be considered as a single critical section.
- Partitions are scheduled by using a predefined scheduling cyclic plan.
- 505 • Inter-partition communication is done through messages.
- There are two kind of virtual communication devices: sampling ports and queuing ports.
- All hypercall services are non-blocking.
- Regarding the capabilities of the partitions, XtratuM defines two kind of partitions: system and standard.
- 510 • Only system partitions are allowed to control the state of the system and other partitions, and to query about them.
- XtratuM is configured off-line and no dynamic objects can be added at run-time.
- The XtratuM configuration file (XM_CF) describes the resources that are allowed to be used by each partition.
- 515 • XtratuM provides a fine grain error detection and a coarse grain fault management.
- It is possible to implement advanced fault analysis techniques in system partitions.
- An I/O Server partition can handle a set of devices used by several partitions.
- XtratuM implements a highly configurable health monitoring and handling system.
- The logs reported by the health monitoring system can be retrieved and analysed by a system partition online.
- 520 • XtratuM provides a tracing service that can be used to both debug partitions and online monitoring.
- The same tracing mechanism is used to handle partition and XtratuM traces.

Chapter 3

Developing Process Overview

XtratuM is a layer of software that extends the capabilities of the native hardware. There are important differences between a classical system and an hypervisor based one. This chapter provides an overview of the XtratuM developing environment. 525

The simplest scenario is composed of two actors: the *integrator* and two *partition developer* or partition supplier. There shall be only one integrator team and one or more partition developer teams (in what follows, we will use “integrator” and “partition developer” for short).

The tasks to be done by the **integrator** are: 530

1. Configure the XtratuM source code (jointly with the resident software). Customise it for the target board (processor model, etc.) and a miscellaneous set of code options and limits (debugging, identifiers length, etc.). See section 8.1 for a detailed description.
2. Build XtratuM: hypervisor binary, user libraries and tools.
3. Distribute the resulting binaries to the partition developers. All partition developers shall use the same binary version of XtratuM. 535
4. Allocate the available system resources to the partitions, according to the resources required to execute each partition:
 - memory areas where each partition will be executed or can use,
 - design the scheduling plan, 540
 - communication ports between partitions,
 - the virtual devices and physical peripherals allocated to each partition,
 - configure the health monitoring,
 - etc.

By creating the `XM.CF` configuration file¹. See section 8.3 for a detailed description. 545

5. Gather the partition images and customisation files from partition developers.
6. Pack all the files (resident software, XtratuM binary, partitions, and configuration files) into the final system image.

The **partition developer** activity:

1. Define the resources required by its application, and send it to the integrator. 550

¹Although it is not mandatory to name “`XM.CF`” the configuration file, we will use this name in what follows for simplicity.

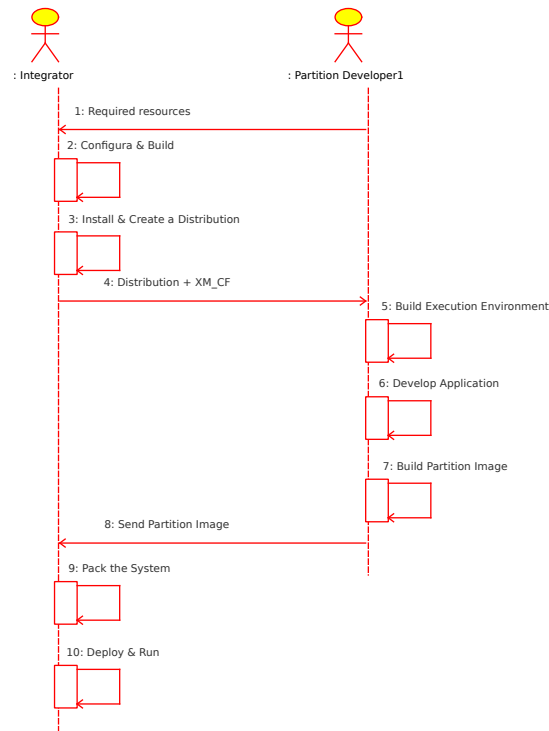


Figure 3.1: Integrator and partition developer interactions.

2. Prepare the development environment. Install the binary distribution created by the integrator.
3. Develop the partition application, according to the system resources agreed by the integrator.
4. Deliver to the integrator the resulting partition image and the required customisation files (if any).

555 There should be an agreement between the integrator and the partition developers on the resources allocated to each partition. The binaries, jointly with the XM_CF configuration file defines the partitioned system. **All partition developers shall use exactly the same XtratuM binaries and configuration files during the development.** Any change on the configuration shall be agreed with the integrator.

560 Since the development of the partitions may be carried out in parallel (or due to intellectual property restrictions), the binary image of some partitions may not be available to a partition developer team. In this case, it is advisable to use dummy partitions to replace those non-available, rather than changing the configuration file.

3.1 Development at a glance

- ① The first step is to build the hypervisor binaries. The integrator shall configure and compile the XtratuM sources to produce:

565 **xm_core.xef:** The hypervisor image which implements the support for partition execution.

libxm.a: A helper library which provides a “C” interface to the para-virtualised services via the hypercall mechanism.

xmc.xsd: The XML schema specification to be used in the XM_CF configuration file.

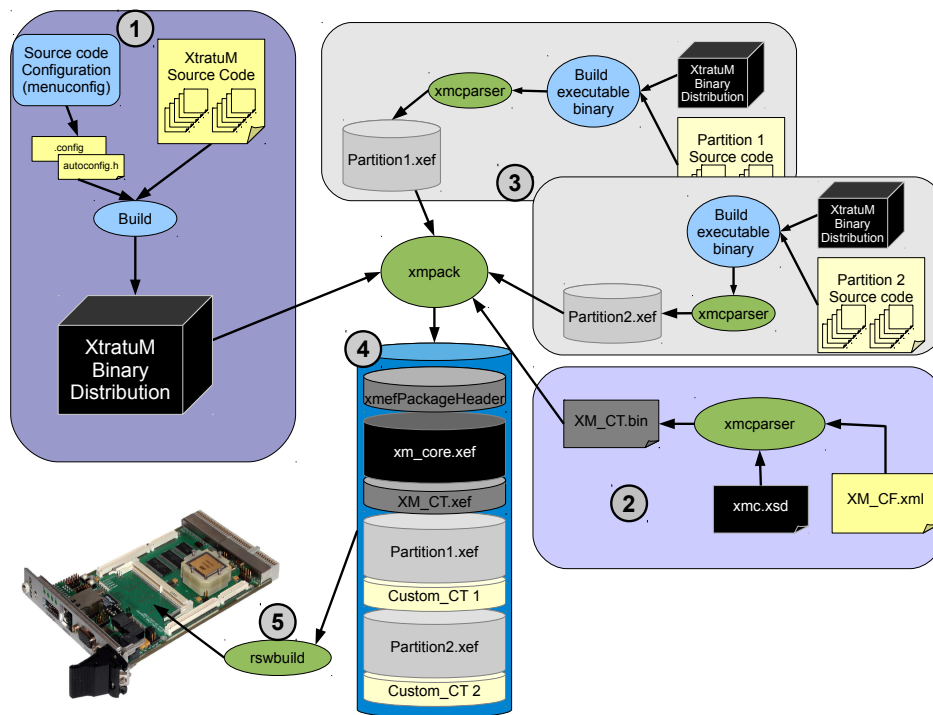


Figure 3.2: The big picture of building a XtratuM system.

tools: A set of tools to manage the partition images and the XM_CF file.

The result of the build process can be prepared to be delivered to the partition developers as a binary distribution. 570

- ② The next step is to define the hypervisor system and resources allocated to each partition. This is done by creating the configuration file XM_CF file.
- ③ Using the binaries resulted from the compilation of XtratuM and the system configuration file, partition developers can implement and test its own partition code by their own. 575
- ④ The tool xmpack is used to build the complete system (hypervisor plus partitions code). The result is a single file called *container*. Partition developers shall replace the image of non-available partitions by a dummy partition. Up to three, customisation files can be attached to each partition.
- ⑤ The container shall be loaded in the target system using the corresponding resident software (or boot loader). For convenience, a resident software is provided. 580

3.2 Building XtratuM

In the first stage, **XtratuM shall be tailored to the hardware available on the board, and the expected workload.** This configuration parameters will be used in the compilation of the XtratuM code to produce a compact and efficient XtratuM executable image. Parameters like the processor model or the memory layout of the board are configured here (see section 8.1).

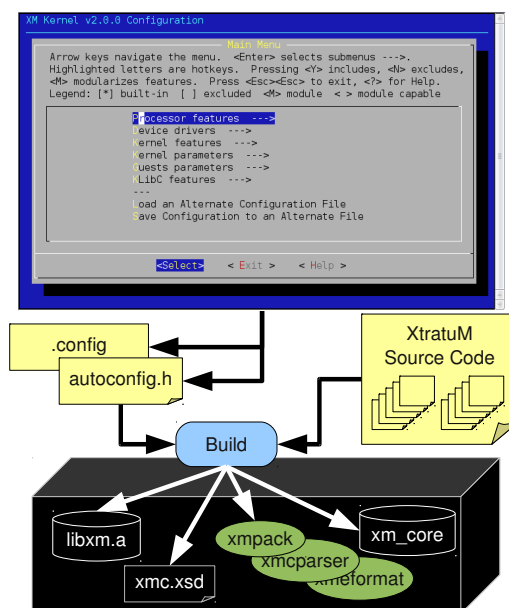


Figure 3.3: Menuconfig process.

585 The configuration interface is the same than the one known as “*menuconfig*” used in the Linux kernel, see figure 3.3. It is a ncurses-based graphic interface to edit the configuration options. The selected choices are stored in two files: a “C” include file named “`core/include/autoconf.h`”; and a makefile include file named “`core/.config`”. Both files contain the same information but with different syntax to be use “C” programs and the in Makefiles respectively.

590 Although it is possible to edit these configuration files, with a plain text editor, it is advisable not to do so; since both files shall be synchronized.

Once configured, the next step is to build XtratuM binaries, which is done calling the command `make`.

Ideally, configuring and compiling XtratuM should be done at the initial phases of the design and should not be changed later.

595 The build process leaves the objects and executables files in the sources directory. Although it is possible to use these files directly to develop partitions it is advisable to install the binaries in a separate read-only directory to avoid accidental modifications of the code. It is also possible to build a TGZ² package with all the files to develop with XtratuM, which can be delivered to the partition developers. See chapter 4.

3.3 System configuration

600 The integrator, jointly with the partition developers, have to define the resources allocated to each partition, by creating the `XM_CF` file. It is an XML file which shall be a valid XML against the XMLSchema defined in section 8.3. Figure 3.4 shows a graphical view of the configuration schema.

The main information contained in the `XM.CF` file is:

605 **Memory:** The amount of physical memory available in the board and the memory allocated to each partition.

²TGZ: Tar Gzipped archive.


Processor: How the processor is allocated to each partition: the scheduling plan.

Peripherals: Those peripherals not managed by XtratuM can be used by one partition. The I/O port ranges and the interrupt line if any.

Health monitoring: How the detected errors are managed by the partition and XtratuM: direct action, delivered to the offending partition, create a log entry, reset, etc. 610

Inter-partition communication: The ports that each partition can use and the channels that link the source and destination ports.

Tracing: Where to store trace messages and what messages shall be traced.

Since XM_CF defines the resources allocated to each partition, this file represents a **contract between the integrator and the partition developers**. A partner (the integrator or any of the partition developers) should not change the contents of the configuration file on its own. All the partners should be aware of the changes and should agree in the new configuration in order to avoid problems later during the integration phase. 615 

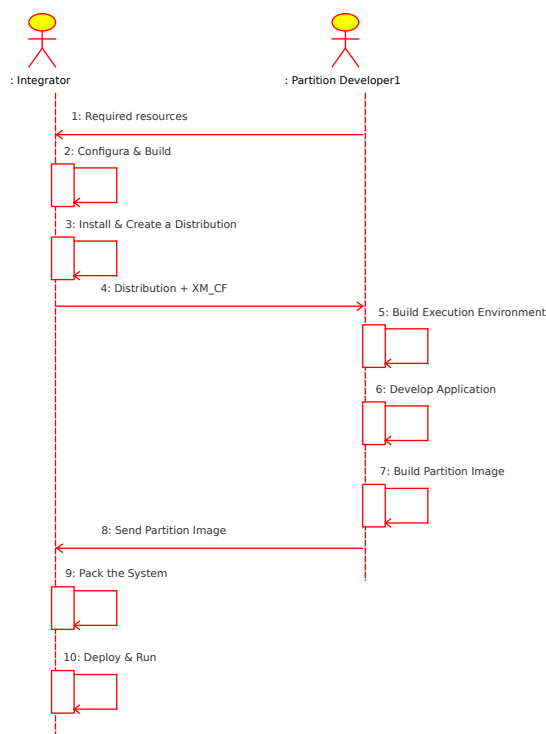


Figure 3.4: Graphical representation of an XML configuration file.

In order to reduce the complexity of the XtratuM hypervisor, the XM_CF is parsed and translated into a binary format which can be directly used by XtratuM. The XML data is translated into a set of initialised data structures ready to be used by XtratuM. Otherwise, XtratuM will need to contain an XML parser to read the XM_CF information. See section /refxmcparser. 620

The resulting configuration binary will be passed to XtratuM as a “customisation” file.

3.4 Compiling partition code

Partition developers should use the XtratuM user library (named `libxm.a` which has been generated during the compilation of the XtratuM source code) to access the para-virtualised services. The resulting 625

binary image of the partition shall be self-contained, that is, it shall not contain linking information. The ABI of the partition binary is described in section 6.

In order to be able to run the partition application, each partition developer require the following files:

630 **libxm.a**: Para-virtualised services. The include files are distributed jointly with the library, and they should be provided by the integrator.

XM.CF.xml: The system configuration file. This file describes the whole system. The same file should be used by all the partners.

635 **xm.core.bin**: The hypervisor executable. This file is also produced by the integrator, and delivered to the other partners.

xmpack: The tool that packs together, into a single system image container, all the *components*.

xmeformat: For converting an ELF file into an XEF one.

xmcparser: The tool to translate the configuration file (XM.CF.xml) into a “C” file which could be compiled to produce the configuration table (XM.CT).

640 Partition developer should use an execution environment as close as possible to the final system: the same processor board and the same hypervisor framework. To achieve this goal, they should use the same configuration file than the one used by the integrator. But the code of other partitions may be replaced by dummy partitions. This dummy partition code executes just a busy loop to waste time.

3.5 Passing parameters to the partitions: customisation files

645 User data can be passed to each partition at boot time. This information is passed to the partition via the *customisation* files.

It is possible to attach up to three customisation files for partition. The content of each customisation file is copied into the partition memory space at boot time (before the partition boots). The buffer where each customisation file is loaded is specified in the partition header. See section 6.

This is the mechanism used by XtratuM to get the compiled XML system configuration.

3.6 Building the final system image

650 In order to ensure that each partition does not depend on, or affects other partitions or the hypervisor, due to shared symbols. The partition binary is not an ELF file. It is a custom format file (called *XEF*) which contains the machine code and the initialized data. See section 6.

655 The *container* is a **single file** which contains all the code, data and configuration information that will be loaded in the target board. In the context of the container, a *component* refers to the set of files that are part of an execution unit (which can be a partition or the hypervisor itself). **xmpack** is a program that reads all the executable images (XEF files) and the configuration/customisation files and produces the container.

660 The container is not a bootable code. That is, it is like a “tar” file which contains a set of files. In order to be able to start the partitioned system, a boot loader shall load the content of the container into the corresponding partition addresses. The utility **rsbuild** creates an bootable ELF file with the resident software and the container.

Chapter 4

Building XtratuM

4.1 Developing environment

XtratuM has been compiled and tested with the following package versions:

Package	Version	Linux package name		Purpose
host gcc	4.2.3 / 4.4.3	gcc-4.2	req	Build host utilities
make	3.81-3build1	make	req	Core
libncurses	5.7+20100313-5	libncurses5-dev	req	Configure source code
binutils	2.18.1build1 / 2.20.1	binutils	req	Core
sparc-toolchain	linux-3.4.4		req	Core
libxml2	2.6.27 / 2.7.6	libxml2-dev	req	Configuration parser
tsim	2.0.10 / 2.0.15		opt	Simulated run
grmon	1.1.31		opt	Deploy and run
perl	5.8.8-12 / 5.10.1		opt	Testing
makeself	2.1.5	makeself	opt	Build self extracting distribution

Packages marked as “req” are required to compile XtratuM. Those packages marked as “opt” are needed to compile or use it in some cases.

665

4.2 Compile XtratuM Hypervisor

It is not required to be supervisor (root) to compile and run XtratuM.

The first step is to prepare the system to compile XtratuM hypervisor.

1. Check that the GNU LIBC Linux GCC 3.4.4 toolchain for SPARC LEON is installed in the system. It can be downloaded from: <http://www.gaisler.com> sparc-linux-3.4.4-x.x.x.tar.bz2
2. Make a deep clean to be sure that there is not previous configurations:

```
$ make distclean
```

3. In the root directory of XtratuM, copy the file `xmconfig.sparc` into `xmconfig`, and edit it to meet your system paths. The variable `XTRATUM_PATH` shall contain the root directory of XtratuM. Also, if the `sparc-linux` toolchain directory is not in the `PATH` then the variable `TARGET_CCPREFIX` shall

670

contain the path to the actual location of the corresponding tools. The prefix of the SPARC v8 tool chain binaries, shall be “sparc-linux-”; otherwise, edit also the appropriate variables.

In the seldom case that the host toolchain is not in the PATH, then it shall be specified in the HOST_CCPREFIX variable.

```
$ cp xmconfig.sparc xmconfig
$ vi xmconfig .....
```

- 675 4. Configure the XtratuM sources. The ncurses5 library is required to compile the configuration tool. In a Debian system with internet connection, the required library can be installed with the following command: `sudo apt-get install libncurses5-dev`.

The configuration utility is executed (compiled and executed) with the next command:

```
$ make menuconfig
```

Note: The menuconfig target configures the XtratuM source code and the resident software. Therefore, two different configuration menus are presented, see section 8.1.

- 680 For running XtratuM in the simulator, select the appropriate processor model from the menuconfig menus.

5. Compile XtratuM sources:

```
$ make
> Configuring and building the "XtratuM hypervisor"
> Building XM Core
  - kernel/sparcv8
  - kernel/mmu
  - kernel
  - klibc
  - klibc/sparcv8
  - objects
  - drivers
> Linking XM Core
   text  data  bss   dec   hex filename
 61949   152  81112 143213 22f6d xm_core
de1daa3d40e5a2171d0b16c6fbbf3c6e xm_core.xef
> Done

> Configuring and building the "User utilities"
> Building XM user
  - libxm
  - libxef
  - tools
  - tools/xmpack
  - tools/xmcparser
  - tools/xmgcov
  - tools/xmbuildinfo
  - tools/rswbuild
  - tools/xmcbuild
  - tools/xef
  - xal
  - bootloaders/rsw
> Done
```


4.3 Generating binary a distribution

The generated files from the compilation process are in source code directories. In order to distribute the compiled binary version of XtratuM to the partition developers, a distribution package shall be generated. There are two distribution formats:

Tar file: It is a compressed tar file with all the XtratuM files and an installation script.

```
$ make distro-tar
```

Self-extracting installer: It is a single executable file which contains the distribution and the installation script.

```
$ make distro-run
```

The final installation is exactly the same regarding the distribution format used.

```
$ make distro-run
.....
> Installing XM in "/tmp/xm-distro.21476/xtratum-3.1.2/xm"
  - Generating XM sha1sums
  - Installing XAL
  - Generating XAL sha1sums
  - Installing XM examples
  - Generating XM examples sha1sums
  - Setting read-only (og-w) permission.
> Done

> Generating XM distribution "xtratum-3.1.2.tar.bz2"
> Done

> Generating self extracting binary distribution "xtratum-3.1.2.run"
> Done
```

685

The files `xtratum-x.x.x.tar.bz2` or `xtratum-x.x.x.run` contains all the files requires to work (develop and run) with the partitioned system. This tar file contains two root directories: `xa1` and `xm`, and an installation script.

The directory `xm` contains the XtratuM kernel and the associated developer utilities. Xal stands for *XtratuM Abstraction Layer*, and contains the partition code to setup a basic “C” execution environment. Xal is provided for convenience, and it is not mandatory to use it. Xal is only useful for those partitions with no operating system.

690

Although XtratuM core and related libraries are compiled for the LEON3 processor, some of the host configuration and deploying tools (`xmcparser`, `xmpack` and `xmeformat`) are host executables. If the computer where XtratuM was compiled on and the computer where it is being installed are different processor architectures (32bit and 64bit), the tools may not run properly.

695



4.4 Installing a binary distribution

Decompress the `xtratum-x.x.x.tar.bz2` file in a temporal directory, and execute the install script. Alternatively, if the distributed file is `xtratum-x.x.x.run` then just execute it.

The install script requires only two parameters:

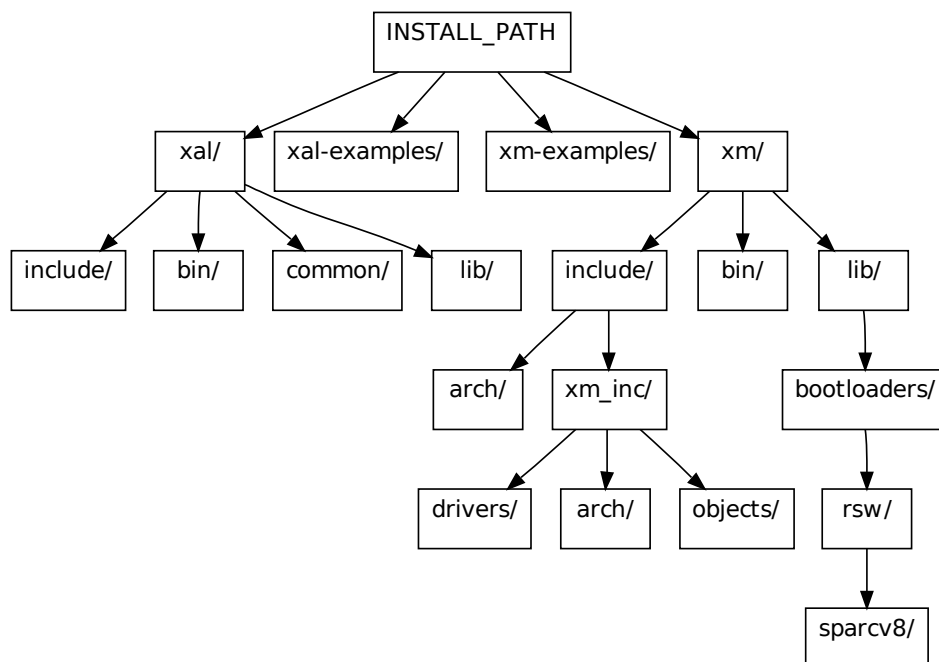


Figure 4.1: Content of the XtratuM distribution.

- 700
1. The installation path.
 2. The path to the SPARC v8 toolchain.

Note that it is assumed that the host toolchain binaries can be located in the PATH variable. It is necessary to provide again the path to the LEON3 toolchain because it may be located in a different place than in the system where XtratuM was build. In any case, it shall be the same version, than the one used to compile XtratuM.

705

```

$ ./xtratum-3.1.2.run
Verifying archive integrity... All good.
Uncompressing XtratuM binary distribution 3.1.2:.....

Starting installation.
Installation log in: /tmp/xtratum-installer.log

1. Select the directory where XtratuM will be installed. The installation
   directory shall not exist.

2. Select the target compiler toolchain binary directory (arch sparc).
   The toolchain shall contain the executables named sparc-linux-*.

3. Confirm the installation settings.

Important: you need write permission in the path of the installation directory.

Continue with the installation [Y/n]? Y

Press [Enter] for the default value or enter a new one.
Press [TAB] to complete directory names.

1.- Installation directory [/opt]: /home/xmuser/xm2env
2.- Path to the sparc toolchain [/opt/sparc-linux-3.4.4/bin/]: /opt/sparc-linux/bin

Confirm the Installation settings:
Selected installation path : /home/xmuser/xm2env
Selected toolchain path : /opt/sparc-linux/bin
  
```

```
3.- Perform the installation using the above settings [Y/n]? Y
Installation completed.
```

Listing 4.1: Output of the self-executable distribution file.

4.5 Compile the Hello World! partition

1. Change to the `INSTALL_PATH/xm-examples/hello_world` directory.
2. Compile the partition:

```
$ make
.....
Created by "iripoll" on "fredes" at "Wed Feb 3 13:21:02 CET 2010"
XM path: " /home/xmuser/xm2env/xm"

XtratuM Core:
  Version: "3.1.2"
  Arch:    "sparcv8"
  File:    "/home/xmuser/xm2env/xm/lib/xm_core.xef"
  Sha1:    "962b930e8278df873599e6b8bc4fcb939eb92a19"
  Changed: "2010-02-03 13:19:51.000000000 +0100"

XtratuM Library:
  Version: "3.1.2"
  File:    "/home/xmuser/xm2env/xm/lib/libxm.a"
  Sha1:    "46f64cf2510646833a320e1e4a8ce20e4cd4e0a9"
  Changed: "2010-02-03 13:19:51.000000000 +0100"

XtratuM Tools:
  File:    "/home/xmuser/xm2env/xm/bin/xmcparser"
  Sha1:    "8fa5dea03739cbb3436d24d5bc0e33b20906c47a"
```

Note that the compilation is quite verbose: the compilation commands, messages, detailed information about the tools libraries used, etc. are printed.

The result from the compilation is a file called “resident_sw”.

3. To run this file just load it in the `tsim` or `grmon` and run (go) it:

```
$ tsim-leon3 -mmu
...
serial port A on stdin/stdout
allocated 4096 K RAM memory, in 1 bank(s)
allocated 16 M SDRAM memory, in 1 bank
allocated 2048 K ROM memory
icache: 1 * 4 kbytes, 16 bytes/line (4 kbytes total)
dcache: 1 * 4 kbytes, 16 bytes/line (4 kbytes total)
tsim> load resident_sw
section: .text, addr: 0x40200000, size 14340 bytes
section: .rodata, addr: 0x40203808, size 790 bytes
section: .container, addr: 0x40203b20, size 48888 bytes
section: .got, addr: 0x4020fa18, size 8 bytes
```

710

```

section: .eh_frame, addr: 0x4020fa20, size 64 bytes
read 38 symbols
tsim> go
resuming at 0x40200d24
XM Hypervisor (3.1 r2)
Detected 50.0MHz processor.
>> HwClocks [LEON clock (1000Khz)]
[CPU:0] >> HwTimer [LEON timer (1000Khz)]
[CPU:0] [sched] using cyclic scheduler
2 Partition(s) created
P0 ("Partition1":0) flags: [ SYSTEM BOOT (0x40080000) ]:
  [0x40080000 - 0x400fffff]
P1 ("Partition2":1) flags: [ SYSTEM BOOT (0x40100000) ]:
  [0x40100000 - 0x4017ffff]
Jumping to partition at 0x40082000
Jumping to partition at 0x40102000
I am Partition2
Hello World!
I am [CPU:0] [HYPERCALL] (0x1) Halted
Partition1
Hello World!
[CPU:0] [HYPERCALL] (0x0) Halted

```

4.6 XtratuM directory tree

```

sources_sparcv8
|-- core/
|   |-- drivers/
|   |-- include/
715 |   |-- kernel/
|   |-- klibc/
|   |-- objects/
|   |-- build.info
|   |-- Kconfig.ver
720 |   |-- Makefile
|   |-- rules.mk
|   |-- xm_core
|   \-- xm_core.xef
|-- docs/
725 |   \-- COPYING
|-- scripts/
|   |-- kconfig/
|   |-- asm-offsets.c
|   |-- asm-offsets.sh*
730 |   |-- errataack.pl*
|   |-- extractinfo*
|   |-- extractinfo.c
|   \-- gencomp.pl*
-- user/
735 |   |-- bin/
|   |-- bootloaders/
|   |-- examples/
|   |-- libxref/
|   |-- libxm/
740 |   |-- tools/
|   |-- xal/
|   |-- install.mk
|   |-- Makefile
|   \-- rules.mk
745 |-- Changelog
|-- config.mk
|-- Makefile
|-- README
|-- version
750 |-- xmconfig
\-- xmconfig.sparc

```

Chapter 5

Partition Programming

This chapter explains how to build a XtratuM partition: partition developer tutorial.

5.1 Implementation requirements

Below is a checklist of what the partition developer and the integrator should take into account when using XtratuM. It is advisable to revisit this list to avoid incorrect assumptions.

Development host: If the computer where XtratuM was compiled on and the computer where it is being installed are different processor architectures (32bit and 64bit), the tools may not run properly. 755

Check that the executable files in `xm/bin` are compatible with the host architecture.

Para-virtualised services: Partition's code shall use the para-virtualised services. The use of native services is considered an error and the corresponding error will be raised. 760

PIT and PCT: In the case of corrupting the Partition Control Table, the result on the faulting partition is undefined. The rest of the partitions are not affected.

Store ordering: XtratuM has been implemented considering that the LEON2 processor is operating in TSO (Total Store Ordering) mode. This is the standard SPARC v8 working mode. If changed to PSO (Partial Store Ordering) mode then random errors will happen. 765

Memory allocation:

- Care shall be taken to avoid overlapping the memory allocated to each partition.
- If MMU is not used, then the partition code shall be linked to work on the allocated memory areas. If the memory allocated in the `XM_CF` file is changed, then the linker script of the affected partition shall be updated accordingly. 770

Reserved names: The prefix “xm”, both in upper and lower case, is reserved for XtratuM identifiers.

Stack management: XtratuM manages automatically the register window of the partitions. The partition code is responsible of initialising the stack pointer to a valid memory area, and reserve enough space to accommodate all the data that will be stored in the stack. Otherwise, a stack overflow may occur. 775

Data Alignment: By default, all data structures passed to or shared with XtratuM shall be aligned to 8 bytes.

Units definition and abbreviations:

“KB” (KByte or Kbyte) is equal to 1024 (2^{10}) bytes.

780 “Kb” (Kbit) is equal to 1024 (2^{10}) bits.

“MB” (MByte or Mbyte) is equal to 1048576 ($1024 \cdot 1024 = 2^{20}$) bytes.

“Mb” (Mbit) is equal to 1048576 ($1024 \cdot 1024 = 2^{20}$) bits.

“Khz” (Kilo Hertz) is equal to 1000 hertz.

“Mhz” (Mega Hertz) is equal to 1000.000 hertz.

785 **XtratuM memory footprint:** XtratuM does not use dynamic memory allocation. Therefore, all internal data structures are declared statically. The size of these data structures are defined during the source code configuration process.

The following configuration parameters are the ones that have an impact on the memory needed by XtratuM:

790 **Maximum identifier length:** Defines the space reserved to store the names of the partitions, ports, scheduling slots and channels.



Kernel stack size: For each partition, XtratuM reserves a kernel stack. Do not reduce the value of this parameter unless you know the implications.

795 **Partition memory areas (if the WPR is used) :** Due to the hardware device (WPR) used to force memory protection, the area of memory allocated to the partitions shall fulfil the next conditions:

- The size shall be greater than or equal to 32KB.
- The size shall be a power of two.
- The start address shall be a multiple of the size.

800 **Configuration of the resident software (RSW):** The information contained in the XM_CF regarding the RSW is not used to configure the RSW itself. That information is used:

- by XtratuM to perform a system cold reset,
- and by the `xmcparser` to check for memory overlaps.

805 **Partition declaration order:** The partition elements, in the XM_CF file, shall be ordered by “id”, and the id’s shall be consecutive starting in zero.

5.2 XAL development environment

XAL is a minimal developing environment to create bare “C” applications. It is provided jointly with the XtratuM core. Currently it is only the minimal libraries and scripts to compile and link a “C” application. More features will be added in the future (mathematic lib, etc.).

810 In the previous versions of XtratuM, XAL was included as part of the examples of XtratuM. It has been moved outside the tree of XtratuM to create an independent developer environment.

When XtratuM is installed, the XAL environment is also installed. It is included in the target directory of the installation path.

```
target_directory
|-- xal                # XAL components
|-- xal-examples      # examples of XAL use
|-- xm
```

```
'-- xm-examples
```

Listing 5.1: Installation tree.

The XAL subtree contains the following elements:

```
xal
|-- bin                # utilities
| |-- xpath
| |-- '--- xpathstart
|-- common            # compilation rules
| |-- config.mk
| |-- config.mk.dist
| |-- rules.mk
|-- include           # headers
| |-- arch
| |-- '--- irqs.h
| |-- assert.h
| |-- autoconf.h
| |-- config.h
| |-- ctype.h
| |-- irqs.h
| |-- limits.h
| |-- stdarg.h
| |-- stddef.h
| |-- stdio.h
| |-- stdlib.h
| |-- string.h
| |-- '--- xal.h
|-- lib                # libraries
| |-- libxal.a
| |-- '--- loader.lds
'--- sha1sum.txt
```

Listing 5.2: XAL subtree.

A XAL partition can:

- Be specified as "system" or "user".
- Use all the XtratuM hypercalls according to the type of partition.
- Use the standard input/output "C" functions: printf, sprintf, etc. The available functions are defined in the include/stdio.h.
- Define interrupt handlers and all services provided by XtratuM.

815

An example of a XAL partition is

```
#include <xm.h>
#include <stdio.h>

#define LIMIT 100

void SpentTime(int n) {
    int i,j;
    int x,y = 1;
    for (i= 0; i <=n; i++) {
        for (j= 0; j <=n; j++) {
```

820

```

        x = x + x - y;
    }
}

void PartitionMain(void) {
    long counter=0;

    printf("[P%d] XAL Partition \n",XM_PARTITION_SELF);
    counter=1;
    while(1) {
        counter++;
        SpentTime(2000);
        printf("[P%d] Counter %d \n",XM_PARTITION_SELF, counter);
    }
    XM_halt_partition(XM_PARTITION_SELF);
}

```

Listing 5.3: XAL partition example.

In the `xal-examples` subtree, the reader can find several examples of XAL partitions and how these examples can be compiled. Next is shown the Makefile file.

```

# XAL_PATH: path to the XTRATUM directory
XAL_PATH=/...../xal

# XMLCF: path to the XML configuration file
XMLCF=xm_cf.sparcv8.xml

# PARTITIONS: partition files (xef format) composing the example
PARTITIONS=partition1.xef partition2.xef ....

all: container.bin resident_sw
include $(XAL_PATH)/common/rules.mk

partition1: dummy_xal.o
    $(LD) -o $@ $^ $(LDFLAGS) -Ttext=$(call xpathstart,1,$(XMLCF))
.....

PACK_ARGS=-h $(XMCORE):xm_cf.xef.xmc \
    -p 0:partition1.xef\
    -p 1:partition2.xef\
    ....

container.bin: $(PARTITIONS) xm_cf.xef.xmc
    $(XMPACK) check xm_cf.xef.xmc $(PACK_ARGS)
    $(XMPACK) build $(PACK_ARGS) $@
    @exec echo -en "> Done [container]\n"

```

Listing 5.4: Makefile.

5.3 Partition definition

A partition is an execution environment managed by the hypervisor which uses the virtualised services.

Each partition consists of one or more concurrent processes (implemented by the operating system of

each partition), sharing access to processor resources based upon the requirements of the application. The partition code can be:

825

- An application compiled to be executed on a bare-machine (bare-application).
- A real-time operating system and its applications.
- A general purpose operating system and its applications.

Partitions need to be *virtualised* to be executed on top of XtratuM. For instance, the partitions cannot manage directly the hardware interrupts (enable/disable interrupts) which have to be replaced by hypercalls¹ to ask for the hypervisor to enable/disable the interrupts.

830

Depending on the type of execution environment, the virtualisation implies:

Bare application The application has to be virtualised using the services provided by XtratuM. The application is designed to run directly on the hardware and it has to be aware about it.

835

Operating system application When the application runs on top of a (real-time) operating system, it uses the services provided by the operating system and does not need to be virtualised. But the operating system has to deal with the virtualisation. The operating system has to be virtualised (ported on top of XtratuM).

5.4 The “Hello World” example

Let’s start with a simple code that is not ready to be executed on XtratuM and needs to be adapted.

840

```
void main() {
    int counter =0;

    xprintf(‘Hello World!\n’);
    while(1) {
        counter++;
        counter %= 100000;
    }
}
```

Listing 5.5: Simple example.

The first step is to initialise the virtual execution environment and call the entry point (`PartitionMain` in the examples) of the partition. The following files are provided as an example of how to build the partition image and initialise the virtual machine.

boot.S: The assembly code where the headers and the entry point are defined.

traps.c: Required data structures: PCT and trap handlers.

845

std.c.c, std.c.h: Minimal “C” support as `mencpy`, `xprintf`, etc.

loader.lds: The linker script that arranges the sections to build the partition image layout.

The `boot.S` file:

¹para-virtualised operations provided by the hypervisor

```

1  #include <xm.h>
2  #include <xm_inc/arch/asm_offsets.h>
3  // #include <xm_inc/hypercalls.h>
4
5  #define STACK_SIZE 8192
6  #define MIN_STACK_FRAME 0x60
7
8  .text
9  .align 8
10
11 .global start, _start
12
13 _start:
14 start:
15     /* Zero out our BSS section. */
16     set _sbss, %o0
17     set _ebss, %o1
18
19 1:
20     st %g0, [%o0]
21     subcc %o0, %o1, %g0
22     bl 1b
23     add %o0, 0x4, %o0
24
25     set __stack_top, %fp
26
27     mov %g1, %o0
28     call init_libxm
29     sub %fp, MIN_STACK_FRAME, %sp
30     ! Set up TBR
31     set write_register32_nr, %o0
32     mov %g0, %o1
33     set _traptab, %o2
34     __XM_HC
35
36     call PartitionMain
37     sub %fp, MIN_STACK_FRAME, %sp
38
39     set 0xFFFFF000, %o0
40     ld [%o0 + 0x1c], %o1
41     mov halt_partition_nr, %o0
42     __XM_HC
43
44 1:    b 1b
45     nop
46
47 ExceptionHandlerAsm:
48     mov sparcv8_get_psr_nr, %o0
49     __XM_AHC
50     mov %o0, %l0
51 !    set sparcv8_flush_regwin_nr, %o0
52 !    __XM_AHC
53     sub %fp, 48, %fp
54     std %l0, [%fp+40]
55     std %l2, [%fp+32]
56     std %g6, [%fp+24]
57     std %g4, [%fp+16]
58     std %g2, [%fp+8]
59     st %g1, [%fp+4]
60     rd %y, %g5
61     st %g5, [%fp]
62
63     mov %l5, %o0
64     call ExceptionHandler
65     sub %fp, 0x80, %sp
66     ld [%fp], %g1
67     wr %g1, %y
68     ld [%fp+4], %g1
69
70     ldd [%fp+8], %g2
71     ldd [%fp+16], %g4
72     ldd [%fp+24], %g6
73     ldd [%fp+32], %l2
74     ldd [%fp+40], %l0
75     add %fp, 48, %fp
76     mov %l0, %o1
77     mov sparcv8_set_psr_nr, %o0
78     __XM_AHC
79     set sparcv8_iret_nr, %o0
80     __XM_AHC
81 ExtIrqHandlerAsm:
82     mov sparcv8_get_psr_nr, %o0
83     __XM_AHC
84     mov %o0, %l0
85 !    set sparcv8_flush_regwin_nr, %o0
86 !    __XM_AHC
87     sub %fp, 48, %fp
88     std %l0, [%fp+40]
89     std %l2, [%fp+32]
90     std %g6, [%fp+24]
91     std %g4, [%fp+16]
92     std %g2, [%fp+8]
93     st %g1, [%fp+4]
94     rd %y, %g5
95     st %g5, [%fp]
96     mov %l5, %o0
97     call ExtIrqHandler
98     sub %fp, 0x80, %sp
99     ld [%fp], %g1
100    wr %g1, %y
101    ld [%fp+4], %g1
102    ldd [%fp+8], %g2
103    ldd [%fp+16], %g4
104    ldd [%fp+24], %g6
105    ldd [%fp+32], %l2
106    ldd [%fp+40], %l0
107    add %fp, 48, %fp
108    mov %l0, %o1
109    mov sparcv8_set_psr_nr, %o0
110    __XM_AHC
111    set sparcv8_iret_nr, %o0
112    __XM_AHC
113
114 HwIrqHandlerAsm:
115     mov sparcv8_get_psr_nr, %o0
116     __XM_AHC
117     mov %o0, %l0
118 !    set sparcv8_flush_regwin_nr, %o0
119 !    __XM_AHC
120     sub %fp, 48, %fp
121     std %l0, [%fp+40]
122     std %l2, [%fp+32]
123     std %g6, [%fp+24]
124     std %g4, [%fp+16]
125     std %g2, [%fp+8]
126     st %g1, [%fp+4]
127     rd %y, %g5
128     st %g5, [%fp]
129     mov %l5, %o0
130     call HwIrqHandler
131     sub %fp, 0x80, %sp
132
133     ld [%fp], %g1
134     wr %g1, %y
135     ld [%fp+4], %g1
136     ldd [%fp+8], %g2

```

```

137     ldd [%fp+16], %g4
138     ldd [%fp+24], %g6
139     ldd [%fp+32], %l2
140     ldd [%fp+40], %l0
141     add %fp, 48, %fp
142     mov %l0, %o1
143     mov sparcv8_set_psr_nr, %o0
144     __XM_AHC
145     set sparcv8_iret_nr, %o0
146     __XM_AHC
147
148     .data
149     .align 8
150     __stack:
151         .fill (STACK_SIZE/4),4,0
152     __stack_top:
153
154     .previous
155
156
157     #define BUILD_IRQ(irqnr) \
158         sethi %hi(HwIrqHandlerAsm), %l4 ;\
159         jmpl %l4 + %lo(HwIrqHandlerAsm), %g0 ;\
160         mov irqnr, %l5 ;\
161         nop
162
163     #define BUILD_TRAP(trapnr) \
164         sethi %hi(ExceptionHandlerAsm), %l4 ;\
165         jmpl %l4 + %lo(ExceptionHandlerAsm), %
            g0 ;\
166         mov trapnr, %l5 ;\
167         nop ;
168
169     #define BAD_TRAP(trapnr) \
170     1:    b 1b ;\
171         nop ;\
172         nop ;\
173         nop
174
175     #define SOFT_TRAP(trapnr) \
176     1:    b 1b ;\
177         nop ;\
178         nop ;\
179         nop
180
181     #define BUILD_EXTIRQ(trapnr) \
182         sethi %hi(ExtIrqHandlerAsm), %l4 ;\
183         jmpl %l4 + %lo(ExtIrqHandlerAsm), %g0
            ;\
184         mov (trapnr+32), %l5 ;\
185         nop
186
187     .align 4096
188     _traptab:
189     ! + 0x00: reset
190     t_reset: b start
191         nop
192         nop
193         nop
194
195     ! + 0x01: instruction_access_exception
196         BUILD_TRAP(0x1)
197
198     ! + 0x02: illegal_instruction
199         BUILD_TRAP(0x2)
200
201     ! + 0x03: privileged_instruction
202         BUILD_TRAP(0x3)
203
204     ! + 0x04: fp_disabled
205         BUILD_TRAP(0x4)
206
207     ! + 0x05: window_overflow
208         BUILD_TRAP(0x5)
209
210
211     ! .....

```

Listing 5.6: /user/examples/sparcv8/boot.S

The `__xmImageHdr` declares the required image header (see section 6) and one partition header²: `__xmPartitionHdr`.

850

The entry point of the partition (the first instruction executed) is labeled `start`. First off, the `bss` section is zeroed; the stack pointer (`%sp` register) is set to a valid address; the address of the partition header is passed to the `libxm` (`call InitLibxm`); the virtual trap table register is loaded with the direction of `__traptab`; and finally the user routine `PartitionMain` is called. If the main function returns, then an endless loop is executed.

855

The remaining of this file contains the trap handler routines. Note that the assembly routines are only provided as illustrative examples, and **should not be used on production application systems**. These trap routines just jump to “C” code which is located in the file `traps.c`.



```

#include <xm.h>
#include <xm_inc/arch/paging.h>
#include "std_c.h"

extern void start(void);

struct xmImageHdr xmImageHdr __XMIHDR = {

```

²Multiple partition headers can be declared to allocate several processors to a single partition (experimental feature not documented).

```

    .sSignature=XMEF_PARTITION_MAGIC,
    .compilationXmAbiVersion=XM_SET_VERSION(XM_ABI_VERSION, XM_ABI_SUBVERSION, XM_ABI_REVISION
    ),
    .compilationXmApiVersion=XM_SET_VERSION(XM_API_VERSION, XM_API_SUBVERSION, XM_API_REVISION
    ),
    .noCustomFiles=0,
    /*
    .customFileTab={ [0]=(struct xefCustomFile){
        .sAddr=(xmAddress_t)0x40180000,
        .size=0,
    },
    },
    */
    .eSignature=XMEF_PARTITION_MAGIC,
};

void __attribute__((weak)) ExceptionHandler(xm_s32_t trapNr) {
    xprintf("exception 0x%x (%d)\n", trapNr, trapNr);
    //XM_halt_partition(XM_PARTITION_SELF);
}

void __attribute__((weak)) ExtIrqHandler(xm_s32_t trapNr) {
    xprintf("extIrq 0x%x (%d)\n", trapNr, trapNr);
    // XM_halt_partition(XM_PARTITION_SELF);
}

void __attribute__((weak)) HwIrqHandler(xm_s32_t trapNr) {
    xprintf("hwIrq 0x%x (%d)\n", trapNr, trapNr);
    // XM_halt_partition(XM_PARTITION_SELF);
}

```

Listing 5.7: /user/examples/common/traps.c

Note that the “C” trap handler functions are defined as “weak”. Therefore, if these symbols are defined elsewhere, the new declaration will replace this one.

The linker script that arranges all the ELF sections is:

```

/*OUTPUT_FORMAT("binary")*/
OUTPUT_FORMAT("elf32-sparc", "elf32-
    sparc", "elf32-sparc")
OUTPUT_ARCH(sparc)
ENTRY(start)

SECTIONS
{
    _sguest = .;
    .text ALIGN (8) : {
        *(.text.init)
        *(.text)
    }
    .rodata ALIGN (8) : {
        *(.rodata)
        *(.rodata.*)
        *(.rodata.*.*)
    }

    .data ALIGN (8) : {
        _sdata = .;
        *(.data)
        _edata = .;
    }

    .bss ALIGN (8) : {
        _sbss = .;
        *(COMMON)
        *(.bss)
        _ebss = .;
    }

    _eguest = .;

    /DISCARD/ :
    {
        *(.note)
        *(.comment*)
    }
}

```

```
}
}
```

Listing 5.8: /user/xal/lib/loader.lds

The section `.text.ini`, which contains the headers, is located at the beginning of the file (as defined by the ABI). The section `.xm_ctl`, which contains the PCT table, is located at the start of the `bss` section to avoid being zeroed at the startup of the partition. The contents of these tables has been initialized by XtratuM before starting the partition. The symbols `_sguest` and `_eguest` mark the Start and End of the partition image.

860

The ported version of the previous simple code is the following:

```
#include "std_c.h"          /* Helper functions */
#include <xm.h>
void PartitionMain () {    /* ‘C’ code entry point. */
    int counter=0;

    xprintf(‘Hello World!\n’);
    while(1) {
        counter++;
        counter %= 100000;
    }
}
```

Listing 5.9: Ported simple example

Listing 5.10 shows the main compilation steps required to generate the final container file, which contains a complete XtratuM system, of a system of only one partition. The partition is only a single file, called `simple.c`. This example is provided only to illustrate the build process. It is advisable to use some of the Makefiles provided in the `xm-examples` (in the installed tree).

865

```
# --> Compile the partition source code: [simple.c] -> [simple.o]
$ sparc-linux-gcc -Wall -O2 -nostdlib -nostdinc -Dsparcv8 -fno-strict-aliasing \
  -fomit-frame-pointer --include xm_inc/config.h --include xm_inc/arch/arch_types.h \
  -I[...]/libxm/include -DCONFIG_VERSION=2 -DCONFIG_SUBVERSION=1 \
  -DCONFIG_REVISION=3 -g -D_DEBUG_ -c -o simple.o simple.c

# --> Link it with the startup (libexamples.a)
$ sparc-linux-ld -o simple simple.o -n -u start -T[...]/lib/loader.lds -L../lib \
  -L[...]/xm/lib --start-group 'sparc-linux-gcc -print-libgcc-file-name' -lxm -lxef \
  -lexamples --end-group -Ttext=0x40080000

# --> Convert the partition ELF to the XEF format.
$ xmeformat build -c simple -o simple.xef

# --> Compile the configuration file.
$ xmcparser -o xm_cf.bin.xmc xm_cf.sparcv8.xml

# --> Convert the configuration file to the XEF format.
$ xmeformat build -c -m xm_cf.bin.xmc -o xm_cf.xef.xmc

# --> Pack all the XEF files of the system into a single container
$ xmpack build -h [...]/xm/lib/xm_core.xef:xm_cf.xef.xmc -p 0:simple.xef container.bin

# --> Build the final bootable file with the resident sw and the container.
$ rswbuild container.bin resident_sw
```

Listing 5.10: Example of a compilation sequence.

The partition code shall be compiled with with the flags `-nostdlib` and `-nostdinc` to avoid using host specific facilities which are not provided by XtratuM. The bindings between assembly and “C” are done considering that not frame pointer is used: `-fomit-frame-pointer`.

All the object files (`traps.o`, `boot.o` and `simple.o`) are linked together, and the text section is positioned in the direction `0x40050000`. This address shall be the same than the one declared in the `XM.CF` file:

```
<PartitionTable>
  <Partition id="0" name="Partition1" flags="system" console="Uart">
    <PhysicalMemoryAreas>
      <Area start="0x40080000" size="512KB"/>
    </PhysicalMemoryAreas>
    <TemporalRequirements duration="500ms" period="500ms"/>
  </Partition>
  <Partition id="1" name="Partition2" flags="system" console="Uart">
    <PhysicalMemoryAreas>
      <Area start="0x40100000" size="512KB" flags=""/>
    </PhysicalMemoryAreas>
    <TemporalRequirements duration="500ms" period="500ms"/>
    <HwResources>
      <Interrupts lines="4"/>
    </HwResources>
  </Partition>
</PartitionTable>
```

Listing 5.11: `/user/xal/examples/hello_world/xm_cf.sparcv8.xml`

In order to avoid inconsistencies between the memory `@Area` attribute of the configuration and the parameter passed to the linker, the `examples/common/xpath` tool³ can be used, from a Makefile, to extract the information from the configuration file.

```
$ cd user/examples/hello_world
$ ../common/xpath -c -f xm_cf.sparcv8.xml /SystemDescription/PartitionTable/
  Partition[1]/PhysicalMemoryAreas/Area[1]/@start
0x40050000
```

Listing 5.12: Using `xpath` to recover to memory area of the first partition.

The attribute `/SystemDescription/PartitionTable/Partition[1]/PhysicalMemoryAreas/Area[1]/@start` is the `xpath` reference to the attribute which defines the first region of memory allocated to the first partition, which in the example is the place where the partition will be loaded.

5.4.1 Included headers

The include header which contains all the definitions and declarations of the `libxm.a` library is `xm.h`. This file depends (includes) also the next list of files:

```
includes
|-- comm.h
|-- endianess.h
|-- hm.h
|-- hypervisor.h
|-- status.h
|-- trace.h
|-- xm.h
|-- xmhypercalls.h
```

³`xpath` is a small shell script frontend to the `xmllint` utility.

```

'-- sparcv8
   |-- atomic\_ops.h
   |-- hypervisor.h
   '-- xhypercalls.h

```

895

5.5 Partition reset

A partition reset is an unconditional jump to the partition entry point. There are two modes to reset a partition: `XM_WARM_RESET` and `XM_COLD_RESET`.

On a warm reset, the state of the partition is mostly preserved. Only the field `resetCounter` of the PCT is incremented, and the field `resetStatus` is set to the value given on the hypercall (see `XM_partition_reset()`).

900

On a cold reset: the PCT table is rebuild; `resetCounter` field is set to zero; and `resetStatus` set to the value given on the hypercall; the communication ports are closed; the timers are disarmed.

5.6 System reset

There are two different system reset sequences:

Warm reset: XtratuM jumps to its entry point. This is basically a software reset.

905

Cold reset: A hardware reset if forced. (See section 8.3.5).

The set of actions done on a warm system reset are still under development.

5.7 Scheduling

5.7.1 Slot identification

A partition can get information about which is the current slot being executed querying its PCT. This information can be used to synchronise the operation of the partition with the scheduling plan.

The information provided is:

910

Slot duration: The duration of the current slot. The value of the attribute “duration” for the current slot.

Slot number: The slot position in the system plan, starting in zero.

Id value: Each slot in the configuration file has a required attribute, named “id”, which can be used to label each slot with a user defined number.

915

The id field is not interpreted by XtratuM and can be used to mark, for example, the slots at the starts of each period.

5.7.2 Managing scheduling plans

A system partition can request a plan switch at any time using the hypercall `XM_set_plan()`. The system will change to the new plan at the end of the current MAF. If `XM_set_plan()` is called several times before the end of the current plan, then the plan specified in the last call will take effect.

The hypercall `XM_get_plan_status()` returns information about the plans. The `xmPlanStatus_t` contains the following fields:

```
typedef struct {
    xmTime_t switchTime;
    xm_s32_t next;
    xm_s32_t current;
    xm_s32_t prev;
} xmPlanStatus_t;
```

Listing 5.13: `/core/include/objects/status.h`

switchTime: The absolute time of the last plan switch request. After a reset (both warm and cold), the value is set to zero.

current: Identifier of the current plan.

next: The plan identifier that will be active on the next major frame. If no plan switch is going to occur, then the value of `next` is equal to the value of `current`.

prev: The identifier of the plan executed before the current one. After a reset (both warm and cold) the value is set to `(-1)`.

5.8 Console output

XtratuM offers a basic service to print a string on the console. This service is provided through a hypercall.

```
XM_write_console("Partition 1: Start execution\n", 29);
```

Listing 5.14: Simple hypercall invocation.

Additionally to this low level hypercall, some function have been created to facilitate the use of the console by the partitions. These functions are coded in `examples/common/std_c.c`. Some of these functions are: `strlen()`, `print_str()`, `xprintf()` which are similar to the functions provided by a `stdio`.

The use of `xprintf()` is illustrated in the next example:

```
#include <xm.h>
#include "std_c.h" // header of the std_c.h

void PartitionMain () { // partition entry point
    int counter=0;

    while(1) {
        counter++;
        if (!(counter%1000))
            xprintf("%d\n", counter);
    }
}
```

Listing 5.15: Ported dummy code 1

`xprintf()` performs some format management in the function parameters and invokes the hypercall which stores it in a kernel buffer. This buffer can be sent to the serial output or other device. 935

5.9 Inter-partition communication

Partitions can send/receive messages to/from other partitions. The basic mechanisms provided are sampling and queuing ports. The use of sampling ports is detailed in this section.

Ports need to be defined in the system configuration file `XM.CF`. Source and destination ports are connected through channels. Assuming that ports and channel linking the ports are defined in the configuration file, the next partition code shows how to use it. 940

`XM_create_sampling_port()` and `XM_create_queuing_port()` hypercalls return *object descriptors*. A object descriptor is an integer, where the 16 least significant bits are a unique id of the port and the upper bits are reserved for internal use.

In this example `partition_1` writes values in the `port1` whereas `partition_2` read them. 945

```
#include <xm.h>
#include "std_c.h"

#define PORT_NAME "port1"
#define PORT_SIZE 48

void PartitionMain () { // partition entry point
    int counter=0;
    int portDesc;

    portDesc=XM_create_sampling_port(PORT_NAME,
                                     PORT_SIZE,
                                     XM_SOURCE_PORT);

    if ( portDesc < 0 ) {
        xprintf("[%s] cannot be created", PORT_NAME);
        return;
    }

    while(1) {
        counter++;
        if ( !(counter%1000) ){
            XM_write_sampling_message(portDesc,
                                     counter, sizeof(counter));
        }
    }
}
```

Listing 5.16: Partition_1

```
#include <xm.h>
#include "std_c.h"

#define PORT_NAME "port2"
#define PORT_SIZE 48

void PartitionMain () { // partition entry point
    int value;
    int previous = 0;
    int portDesc;
    xm_u32_t flags;

    portDesc=XM_create_sampling_port(PORT_NAME,
                                     PORT_SIZE,
                                     XM_DESTINATION_PORT);

    if ( portDesc < 0 ) {
        xprintf("[%s] cannot be created", PORT_NAME);
        return;
    }

    while(1) {
        XM_read_sampling_message(portDesc,
                                &value,
                                sizeof(value),
                                &flags);

        if (!(value == previous)){
            xprintf("%d\n", value);
            previous = value;
        }
    }
}
```

Listing 5.17: Partition_2

An interesting exercise is to determine which values will be printed.

5.9.1 Message notification

When a message is sent into a queuing port, or written into a sampling port, XtratuM triggers the extended interrupt `XM_VT_EXT_OBJDESC`. By default, this interrupt is masked when the partition boots.

5.10 Peripheral programming

The LEON2 processor implements a memory-mapped I/O for performing hardware input and output operations to the peripherals. There are two hypercalls to access I/O registers: `XM_sparcv8_inport()` and `XM_sparcv8_outport()`.

In order to be able to access (read from or write to) hardware I/O port the corresponding ports have to be allocated to the partition in the `XM_CF` configuration file.

There are two methods to allocate ports to a partition in the configuration file:

Range of ports: A range of I/O ports, with no restriction, allocated to the partition. The `Range` element is used.

Restricted port: A single I/O port with restrictions on the values that the partition is allowed to write in. The `Restricted` element is used in the configuration file. There are two kind of restrictions that can be specified:

Bitmask: Only those bits that are set, can be modified by the partition. In the case of a read operation only those bits set in the mask will be returned to the partition; the rest of the bits will be resetted. Attribute `mask`.

The attribute `mask` is optional. A restricted port declaration with no attribute, is equivalent to declare a range of ports of size one. In the case that both, the bitmap and the range of values, are specified then the bitmap is applied first and then the range is checked.

The implementation of the bit mask is done as follows:

```
oldValue=LoadIoReg(port);
StoreIoReg(port, ((oldValue&~(mask))|(value&mask)));
```

Listing 5.18: `/core/kernel/sparcv8/hypercalls.c`

First off, the port is read, to get the value of the bits not allocated to the partitions, then the bits which have to be modified are changed, and finally the value is written back.



The read operation shall not cause side effects on the associated peripheral. For example, some devices may interpret as interrupt acknowledge to read from a control port. Another source of errors may happen when the restricted is implemented as an open collector output. In this case, if the pin is connected to an external circuit which forces a low voltage, then the value read from the io port is not the same than the value previous written.

The following example declares a range of ports and two restricted ones.

```
<Partition . . . . . >
  <HwResources>
    <IoPorts>
      <Restricted address="0x3000" mask="0xff" />
    </IoPorts>
  </HwResources>
</Partition>
```

If the bitmask restriction is used, then the bits of the port that are not set in the mask can be allocated to other partitions. This way, it is possible to perform a fine grain (bit level) port allocation to partitions. That is a single ports can be safely shared among several partitions.

975

5.11 Traps, interrupts and exceptions

5.11.1 Traps

A partition can not directly manage processor traps. XtratuM provides a para-virtualized trap system called *virtual traps*. XtratuM defines 256+32 traps. The first 256 traps correspond directly with to the hardware traps. The last 32 ones are defined by XtratuM.



The structure of the virtual trap table mimics the native trap table structure. Each entry is a 16 bytes large and contains the trap handler routine (which in the practice, is a branch or jump instruction to the real handler).

980

At boot time (or after a reset) a partition shall setup its virtual trap table and load the virtual \$tbr register with the start address of the table. The \$tbr register can be managed with the hypercalls: XM_read_register32() and XM_write_register32() with the appropriate register name:

```
#define TBR_REG32 0
```

Listing 5.19: /core/include/arch/processor.h

If a trap is delivered to the partition and there is not a valid virtual trap table, then the health monitoring event XM_EM_EV_PARTITION_UNRECOVERABLE is generated.



5.11.2 Interrupts

In order to properly manage a peripheral, a partition can request to manage directly a hardware interrupt line. To do so, the interrupt line shall be allocated to the partition in the configuration file.

985

There are two groups of virtual interrupts:

Hardware interrupts: Correspond to the native hardware interrupts. Note that SPARC v8 defines only 15 interrupts (from 1 to 15), but XtratuM reserves 32 for compatibility with other architectures.

Interrupt 1 to 15 are assigned to traps 0x11 to 0x1F respectively (as in the native hardware).

990

Extended interrupts: Correspond to the XtratuM extended interrupts.

These interrupts are assigned from traps 0xE0 to 0xFF.

```
#define XM_VT_HW_FIRST          (0)
#define XM_VT_HW_LAST          (31)
#define XM_VT_HW_MAX           (32)

#define XM_VT_HW_INTERNAL_BUS_TRAP_NR (1+XM_VT_HW_FIRST)
#define XM_VT_HW_UART2_TRAP_NR      (2+XM_VT_HW_FIRST)
#define XM_VT_HW_UART1_TRAP_NR      (3+XM_VT_HW_FIRST)
#define XM_VT_HW_IO_IRQ0_TRAP_NR    (4+XM_VT_HW_FIRST)
#define XM_VT_HW_IO_IRQ1_TRAP_NR    (5+XM_VT_HW_FIRST)
#define XM_VT_HW_IO_IRQ2_TRAP_NR    (6+XM_VT_HW_FIRST)
#define XM_VT_HW_IO_IRQ3_TRAP_NR    (7+XM_VT_HW_FIRST)
```

```

#define XM_VT_HW_TIMER1_TRAP_NR    (8+XM_VT_HW_FIRST)
#define XM_VT_HW_TIMER2_TRAP_NR    (9+XM_VT_HW_FIRST)
#define XM_VT_HW_DSU_TRAP_NR       (11+XM_VT_HW_FIRST)
#define XM_VT_HW_PCI_TRAP_NR       (14+XM_VT_HW_FIRST)

#define XM_VT_EXT_FIRST            (0)
#define XM_VT_EXT_LAST             (31)
#define XM_VT_EXT_MAX              (32)

#define XM_VT_EXT_HW_TIMER         (0+XM_VT_EXT_FIRST)
#define XM_VT_EXT_EXEC_TIMER       (1+XM_VT_EXT_FIRST)
#define XM_VT_EXT_WATCHDOG_TIMER   (2+XM_VT_EXT_FIRST)
#define XM_VT_EXT_SHUTDOWN         (3+XM_VT_EXT_FIRST)
#define XM_VT_EXT_OBJDESC          (4+XM_VT_EXT_FIRST)

#define XM_VT_EXT_CYCLIC_SLOT_START (8+XM_VT_EXT_FIRST)

#define XM_VT_EXT_MEM_PROTECT      (16+XM_VT_EXT_FIRST)

/* Inter-Partition Virtual Interrupts */
#define XM_MAX_IPVI 8
#define XM_VT_EXT_IPVIO          (24+XM_VT_EXT_FIRST)
#define XM_VT_EXT_IPVI1          (25+XM_VT_EXT_FIRST)
#define XM_VT_EXT_IPVI2          (26+XM_VT_EXT_FIRST)
#define XM_VT_EXT_IPVI3          (27+XM_VT_EXT_FIRST)
#define XM_VT_EXT_IPVI4          (28+XM_VT_EXT_FIRST)
#define XM_VT_EXT_IPVI5          (29+XM_VT_EXT_FIRST)
#define XM_VT_EXT_IPVI6          (30+XM_VT_EXT_FIRST)
#define XM_VT_EXT_IPVI7          (31+XM_VT_EXT_FIRST)

```

Listing 5.20: /core/include/guest.h

Both, extended and hardware interrupts can be routed to a different interrupt vector through the `XM_route_irq()` hypercall, this hypercall enables a partition to select the most suitable vector to be raised.

995

All hardware and extended interrupts can be masked through the following hypercalls: `XM_clear_irqmask()` and `XM_set_irqmask()`. Besides, all these set of interrupts can be globally disabled/enable by using the `XM_sparc_get_psr()` and `XM_sparc_set_psr()` hypercalls, mimicking the way the underneath architecture works.

An example of interrupts:

```

#include "std_c.h"
#include <xm.h>

void IrqHandler(int irqnr) {
    xprintf("Hardware irq: %d\nHalting.....\n", irqnr);
    XM_halt_partition(XM_PARTITION_SELF);
}

void ExtIrqHandler(int irqnr) {
    xprintf("Extended irq: %d\nHalting.....\n", irqnr);
    XM_halt_partition(XM_PARTITION_SELF);
}

void PartitionMain(void) {
    xm_s32_t err;
    xmTime_t oneshoot;

    XM_get_time(XM_HW_CLOCK , &oneshoot);

```

```

oneshoot += (xmTime_t)100000; // 1 second
XM_set_timer(XM_HW_CLOCK , oneshoot, 0);

XM_enable_irqs();
err = XM_unmask_irq(XM_VT_EXT_HW_TIMER);
xprintf("Unmask extirq XM_VT_EXT_HW_TIMER: %d\n", err);
err = XM_unmask_irq( XM_VT_HW_IO_IRQO_TRAP_NR );
xprintf("Unmask hwirq 4: %d\n", err);

while(1) {
    XM_idle_self();
}
}

```

Listing 5.21: /user/examples/hwirq/hwirq.c

5.11.3 Exceptions

Exceptions are the traps triggered by the processor in response to an internal condition. Some exceptions are caused by normal operation of the processor (e.g. register window over/underflow) but others are caused by abnormal situations (e.g. invalid instruction). 1000

Error related exception traps, are managed by XtratuM thorough the health monitoring system.

```

#define DATA_STORE_ERROR 0x2b // 0
#define INSTRUCTION_ACCESS_MMU_MISS 0x3c // 1
#define INSTRUCTION_ACCESS_ERROR 0x21 // 2
#define R_REGISTER_ACCESS_ERROR 0x20 // 3
#define INSTRUCTION_ACCESS_EXCEPTION 0x1 // 4
#define PRIVILEGED_INSTRUCTION 0x03 // 5
#define ILLEGAL_INSTRUCTION 0x2 // 6
#define FP_DISABLED 0x4 // 7
#define CP_DISABLED 0x24 // 8
#define UNIMPLEMENTED_FLUSH 0x25 // 9
#define WATCHPOINT_DETECTED 0xb // 10
// #define WINDOW_OVERFLOW 0x5
// #define WINDOW_UNDERFLOW 0x6
#define MEM_ADDRESS_NOT_ALIGNED 0x7 // 11
#define FP_EXCEPTION 0x8 // 12
#define CP_EXCEPTION 0x28 // 13
#define DATA_ACCESS_ERROR 0x29 // 14
#define DATA_ACCESS_MMU_MISS 0x2c // 15
#define DATA_ACCESS_EXCEPTION 0x9 // 16
#define TAG_OVERFLOW 0xa // 17
#define DIVISION_BY_ZERO 0x2a // 18

```

Listing 5.22: /core/include/arch/irqs.h

If the health monitoring action associated with the HM event is XM_HM_AC_PROPAGATE, then the same trap number is propagated to the partition as a virtual trap. The partition code is then in charge of handling the error. 1005

5.12 Clock and timer services

XtratuM provides the `XM_get_time()` hypercall to read the time from a clock, and the `XM_set_timer()` hypercall to arm a timer.

There are two clocks available:

```
#define XM_HW_CLOCK (0x0)
#define XM_EXEC_CLOCK (0x1)
```

Listing 5.23: `/core/include/hypercalls.h`

XtratuM provides one timer for each clock. The timers can be programmed in one-shot or in periodic mode. Upon expiration, the extended interrupts `XM_VT_EXT_HW_TIMER` and `XM_VT_EXT_EXEC_TIMER` are triggered. These extended interrupts correspond with traps (`256+XM_VT_EXT_HW_TIMER`) and (`256+XM_VT_EXT_EXEC_TIMER`) respectively.

5.12.1 Execution time clock

The clock `XM_EXEC_CLOCK` only advances while the partition is being executed or while XtratuM is executing a hypercall requested by the partition. The execution time clock computes the total time used by the target partition.

This clock relies on the `XM_HW_CLOCK`, and so, its resolution is also $1\mu\text{sec}$. Its precision is not as accurate as that of the `XM_HW_CLOCK` due to the errors introduced by the partition switch.



The execution time clock does not advance when the partition gets idle or suspended. Therefore, the `XM_EXEC_CLOCK` clock should not be used to arm a timer to wake up a partition from an idle state.

The code below computes the temporal cost of a block of code.

```
#include <xm.h>
#include "std_c.h"

void PartitionMain() {
    xmTime_t t1, t2;

    XM_get_time(XM_EXEC_CLOCK, &t1);
    // code to be measured
    XM_get_time(XM_EXEC_CLOCK, &t2);
    xprintf("Initial time: %lld, final time: %lld", t1, t2);
    xprintf("Difference: %lld\n", t2-t1);
    XM_halt_partition(XM_PARTITION_SELF);
}
```

5.13 Processor management

Currently only the `$tbr` processor control register has been virtualised. This register should be loaded (with the hypercall `XM_write_register32()`) with the address of the partition trap table. This operation is usually done only once when the partition boots (see listing 5.6).

5.13.1 Managing stack context

The SPARC v8 architecture (following the RISC ideas) tries to simplify the complexity of the processor by moving complex management tasks to the compiler or the operating system. One of the most particular features of SPARC v8 is the register window concept, and how it should be managed.

Both, register window overflow and underflow cause a trap. This trap has to be managed in supervisor mode and with traps disabled (bit ET in the \$psr register is unset) to avoid overwriting valid registers. It is not possible to emulate efficiently this behaviour. 1025

XtratuM provides a transparent management of the stack. Stack overflow and underflow is directly managed by XtratuM without the intervention of the partition. The partition code shall load the stack register with valid values. 1030

In order to implement a context switch inside a partition (if the partition is a multi-thread environment), the function `XM_sparcv8_flush_regwin()` can be used to flush (spill) all register windows in the current CPU stack. After calling this function, all the register windows, but the current one, are stored in RAM and then marked as free. The partition context switch code should basically carry out the next actions: 1035

1. call `XM_sparcv8_flush_regwin()`
2. store the current “g”, “i” and “l” registers in the stack
3. switch to the new thread’s stack and
4. restore the same set of registers.

Note that there is no complementary function to reload (fill) the registers, because it is done automatically. 1040

The `XM_sparcv8_flush_regwin()` service can also be used set the processor in a know state before executing a block of code. All the register windows will be clean and no window overflow will happen during the next 7 nested function calls.

5.14 Tracing

5.14.1 Trace messages

The hypercall `XM_trace_event()` stores a trace message in the partition’s associated buffer. A trace message is a `xmTraceStatus_t` structure which contains a `opCode` and an associated user defined data: 1045

```
typedef struct {
    xmTraceOpCode_t opCode;
    xm_u32_t reserved;
    xmTime_t timeStamp;
    union {
        xm_u32_t word[4];
        char str[16];
        struct {
            xmId_t partitionId;
            xm_s32_t newPlanId;
        } auditEvent;
    };
} xmTraceEvent_t;
```

Listing 5.24: /core/include/objects/trace.h

The type `xmTraceOpCode_t` is a 32bit value with the following bit fields:

```
typedef struct {
    xm_u32_t code:13, criticality:3, moduleId:8, partitionId:8;
#define XM_TRACE_UNRECOVERABLE 0x3 // This level triggers a health
                                   // monitoring fault
#define XM_TRACE_WARNING 0x2
#define XM_TRACE_DEBUG 0x1
#define XM_TRACE_NOTIFY 0x0
} xmTraceOpCode_t;
```

Listing 5.25: /core/include/objects/trace.h

partitionId: Identify the partition who issued the trace event. This field is automatically filled by XtratuM. The value `XM_HYPERVISOR_ID` is used to identify XtratuM traces.

moduleId: For the traces issued by the partitions, this field is user defined.

1050 For the traces issued by XtratuM, this field identifies an internal subsystem:

TRACE_MODULE_HYPERVISOR Traces related to XtratuM core events.

TRACE_MODULE_PARTITION Traces related to partition operation.

TRACE_MODULE_SCHED Traces concerning scheduling.

code: This field is user defined for the traces issued by the partitions. For the traces issued by XtratuM, the values of the `code` field depends on the value of the `moduleId`:

1055

If moduleId = TRACE_MODULE_HYPERVISOR

TRACE_EV_HYP_HALT: The hypervisor is about to halt.

TRACE_EV_HYP_RESET: The hypervisor is about to perform a software reset.

TRACE_EV_HYP_AUDIT_INIT: The first message after the audit startup.

1060

If moduleId = TRACE_MODULE_PARTITION

The field `auditEvent.partitionId` has the identifier of the affected partition. The recorded events are:

TRACE_EV_PART_SUSPEND: The affected partition has been suspended.

TRACE_EV_PART_RESUME: The affected partition has been resumed.

1065

TRACE_EV_PART_HALT: The affected partition has been halted.

TRACE_EV_PART_SHUTDOWN: A shutdown extended interrupt has been delivered to the partition.

TRACE_EV_PART_IDLE: The affected partition has been set in idle state.

TRACE_EV_PART_RESET: The affected partition has been reset.

1070

If moduleId = TRACE_MODULE_SCHED

The field `auditEvent.partitionId` has the identifier of the partition that requested the plan switch; or `XM_HYPERVISOR_ID` if the plan switch is the consequence of the `XM_HM_AC_SWITCH_TO_MAINTENANCE` health monitoring action.

The field `auditEvent.newPlanId` has the identifier of the new plan.

1075

TRACE_EV_SCHED_CHANGE_REQ: A plan switch has been requested.


TRACE_EV_SCHED_CHANGE_COMP: A plan switch has been carried out.

criticality: Determines the importance/criticality of the event that motivated the trace message. Next are the intended use of the levels:

XM_TRACE_NOTIFY A notification messages of the progress of the application at coarse-grained level. 1080

XM_TRACE_DEBUG A detailed information message intended to be used for debugging during the development phase.

XM_TRACE_WARNING Traces which informs about potentially harmful situations.

XM_TRACE_UNRECOVERABLE Traces of this level are managed also by the health monitoring sub-system: a user HM event is generated, and handled according to the HM configuration. Note that both, the normal partition trace message is stored, and the HM event is generated. 1085 

Jointly with the `opCode` and the user data, the `XM_trace_event()` function has a bitmask parameter that is used to filter out trace events. If the logical AND between the bitmask parameter and the bitmask of the `XM_CF` configuration file is not zero then the trace event is logged; otherwise it is discarded. Traces of `XM_TRACE_UNRECOVERABLE` critically always raises a health monitoring event regarding the bitmask. 1090

5.14.2 Reading traces

Only one system partition can read from a trace stream. A standard partition **can not read its own trace messages**, it is only allowed to store traces on it.

If the trace stream is stored in a buffer (RAM or FLASH). When the buffer is full, the oldest events are overwritten.

5.14.3 Configuration

XtratuM statically allocates a block of memory to store all traces. The amount of memory reserved to store traces is a configuration parameter of the sources (see section 8.1). 1095

In order to be able to store the traces of a partition, as well as the traces generated by XtratuM, it has to be properly configured in the `XM_CF` configuration file. The `bitmask` attribute is used to filter which traces are stored.

```
<Trace device="MemDisk0" bitmask="0x00000005"/>
```

Listing 5.26: `/user/examples/sched_events/xm_cf.sparcv8.xml`

The traces recoded by XtratuM can be selected (masked) at module granularity.

```
#define TRACE_BM_HYPERVISOR 0x1
#define TRACE_BM_PARTITION 0x2
#define TRACE_BM_SCHED 0x4
```

Listing 5.27: `/core/include/objects/trace.h`

In the example of listing 5.26, the `TRACE_BM_HYPERVISOR` and `TRACE_BM_SCHED` events will be recorded but not `TRACE_BM_PARTITION`. 1100

5.15 System and partition status

The hypercalls `XM_get_partition_status()` and `XM_get_system_status()` return information about a given partition and the system respectively.

The data structure returned are:

```
typedef struct {
    /* Current state of the partition: ready, suspended ... */
    xm_u32_t state;
#define XM_STATUS_IDLE 0x0
#define XM_STATUS_READY 0x1
#define XM_STATUS_SUSPENDED 0x2
#define XM_STATUS_HALTED 0x3
    /* Number of virtual interrupts received. */
    xm_u64_t noVirqs; /* [[OPTIONAL]] */
    /* Reset information */
    xm_u32_t resetCounter;
    xm_u32_t resetStatus;
    xmTime_t execClock;
    /* Total number of partition messages: */
    xm_u64_t noSamplingPortMsgsRead; /* [[OPTIONAL]] */
    xm_u64_t noSamplingPortMsgsWritten; /* [[OPTIONAL]] */
    xm_u64_t noQueuingPortMsgsSent; /* [[OPTIONAL]] */
    xm_u64_t noQueuingPortMsgsReceived; /* [[OPTIONAL]] */
} xmPartitionStatus_t;
```

Listing 5.28: Partition status.

```
typedef struct {
    xm_u32_t resetCounter;
    /* Number of HM events emitted. */
    xm_u64_t noHmEvents; /* [[OPTIONAL]] */
    /* Number of HW interrupts received. */
    xm_u64_t noIrqs; /* [[OPTIONAL]] */
    /* Current major cycle interation. */
    xm_u64_t currentMaf; /* [[OPTIONAL]] */
    /* Total number of system messages: */
    xm_u64_t noSamplingPortMsgsRead; /* [[OPTIONAL]] */
    xm_u64_t noSamplingPortMsgsWritten; /* [[OPTIONAL]] */
    xm_u64_t noQueuingPortMsgsSent; /* [[OPTIONAL]] */
    xm_u64_t noQueuingPortMsgsReceived; /* [[OPTIONAL]] */
} xmSystemStatus_t;
```

Listing 5.29: System status.

The field `execClock` of a partition is the execution time clock of the target partition. The rest of the fields are self explained.

Those fields commented as `[[OPTIONAL]]` contain valid data only if XtratuM has been compiled with the flag “Enable system/partition status accounting” enabled.

5.16 Memory management

XtratuM implements a flat memory space on the SPARC v8 architecture (LEON2 and LEON3 processors). The addresses generated by the control unit are directly emitted to the memory controller without any translation. Therefore, **each partition shall be compiled and linked to work on the designated memory range**. The starting address and the size of each partition is specified in the system configuration file.



1110

Two different hardware features can be used to implement memory protection:

Write Protection Registers (WPR): In the case that there is no MMU support, then it is possible to use the WPR device of the LEON2 and LEON3 processors. The WPR device can be programmed to raise a trap when the processor tries to write on a configured address range. 1115

Since read memory operations are not controlled by the WPR, it is not possible to enforce complete (read/write) memory isolation in this case. Also, due to the internal operation of the WPR device, all the memory allocated to each partition has to be contiguous and has to meet the following conditions: 1120

- The size shall be greater than or equal to 32KB.
- The size shall be a power of two.
- The start address shall be a multiple of the size.

Memory Management Unit (MMU): If the processor has MMU, and XtratuM has been compiled to use it, then fine grain (page size) memory protection is provided. In this case one or more areas of memory can be allocated to each partition. 1125

The MMU is used only as a MPU (memory protection unit), i.e., the virtual and physical addresses are the same. Only the protection bits of the pages are used. As a result, each partition shall be compiled and linked to the designated addresses where they will be loaded and executed.

The memory protection mechanism employed is a source code configuration option. See section 8.1. 1130

The memory areas allocated to a partition are defined in the XM_CF file. The executable image shall be linked to be executed in those allocated memory areas.

The `XM_get_physmem_map()` returns the set of memory areas allocated to the partition. **Available since XtratuM 3.1.**

5.17 Releasing the processor

In some situations, a partition is waiting for a new event to execute a task. If no more tasks are pending to be executed, then the partition can become idle. The idle partition becomes ready again when an interrupt is received. 1135

The partition can inform XtratuM about its idle state (see `XM_idle_self()`). In the current implementation, XtratuM does nothing while a partition is idle, that is, other partitions are not executed; but it opens the possibility to use this wasted time in internal bookkeeping or other maintenance activities. Also, energy saver actions can be done during this idle time. 1140

Since XtratuM delivers an event on every new slot, the idle feature can also be used to synchronise the operation of the partition with the scheduling plan.

5.18 Partition customisation files

A partition is composed of a binary image (code and data) and, zero or more additional files (customisation files). To ease the management of these additional files, the header of the partition image (see section 6.4.1) holds the fields `noModules` and `moduleTab`, where the first is the number of additional files which have to be loaded and the second is an array of data structure which defines the loading address and the sizes of these additional files. During its creation, the partition is responsible for filling these fields with the address of a pre-allocated memory area inside its memory space. 1145

1150 These information shall be used by the loader software, for instance the resident software or a manager system partition, in order to know the place where to copy into RAM these additional files. If the size of any of these files is larger than the one specified on the header of the partition or the memory address is invalid, then the loading process shall fail.

1155 These additional files shall be accessible by part of the loader software. For example, they must be packed jointly with the partition binary image by using the `xmpack` tool.

5.19 Assembly programming

This section describes the assembly programming convention, in order to invoke the XtratuM hypercalls.

The register assignment convention for calling a hypercall is:

`%o0` Holds the hypercall number.

`%o1 - %o5` Holds the parameters to the hypercall.

1160 Once the processor registers have been loaded, a `ta` instruction to the appropriate software trap number shall be called, see section 6.2.

The return value is stored in register `%o0`.

For example, following assembly code calls the `XM_get_time(xm_u32_t clockId, xmTime_t *time)`:

```
mov 0xa, %o0; __GET_TIME_NR
mov %i0, %o1
mov %i1, %o2
ta 0xf0
cmp %o0, 0; XM_OK == 0
bne <error>
```

1165 In SPARC v8, the `get_time_nr` constant has the value “0xa”; “%i0” holds the clock id; and “%i1” is a pointer which points to a `xmTime_t` variable. The return value of the hypercall is stored in “%o0” and then checked if `XM_OK`.

Below is the list of normal hypercall number constants (listing 5.30) and assembly hypercalls (listing 5.31):

```
#define __MULTICALL_NR 0
#define __HALT_PARTITION_NR 1
#define __SUSPEND_PARTITION_NR 2
#define __RESUME_PARTITION_NR 3
#define __RESET_PARTITION_NR 4
#define __SHUTDOWN_PARTITION_NR 5
#define __HALT_SYSTEM_NR 6
#define __RESET_SYSTEM_NR 7
#define __IDLE_SELF_NR 8
#define __WRITE_REGISTER32_NR 9
#define __GET_TIME_NR 10
#define __SET_TIMER_NR 11
#define __READ_OBJECT_NR 12
#define __WRITE_OBJECT_NR 13
#define __SEEK_OBJECT_NR 14
#define __CTRL_OBJECT_NR 15

#define __CLEAR_IRQ_MASK_NR 16
#define __SET_IRQ_MASK_NR 17
```

```
#define __FORCE_IRQS_NR 18
#define __CLEAR_IRQS_NR 19
#define __ROUTE_IRQ_NR 20

#define __UPDATE_PAGE32_NR 21
#define __SET_PAGE_TYPE_NR 22
#define __RAISE_IPVI_NR 23

#define sparcv8_atomic_add_nr 24
#define sparcv8_atomic_and_nr 25
#define sparcv8_atomic_or_nr 26
#define sparcv8_inport_nr 27
#define sparcv8_outport_nr 28
```

Listing 5.30: `xm_inc/arch/hypercalls.h`

```
#define sparcv8_iret_nr 0
#define sparcv8_flush_regwin_nr 1
#define sparcv8_get_psr_nr 2
#define sparcv8_set_psr_nr 3
#define sparcv8_flush_cache_nr 4
```

```
#define sparcv8_flush_tlb_nr 5
#define sparcv8_set_pil_nr 6
#define sparcv8_clear_pil_nr 7
```

Listing 5.31: xm.inc/arch/hypercalls.h

The file “core/include/sparcv8/hypercalls.h” has additional services for the SPARC v8 architecture.

5.19.1 The object interface

XtratuM implements internally a kind of virtual file system (as the /dev directory). Most of the libxm hypercalls are implemented using this file system. The hypercalls to access the objects are used internally by the libxm and shall not be used by the programmer. They are listed here just for completeness:

```
extern __stdcall xm_s32_t XM_read_object(xmObjDesc_t objDesc, void *buffer, xm_u32_t size, xm_u32_t *flags);
extern __stdcall xm_s32_t XM_write_object(xmObjDesc_t objDesc, void *buffer, xm_u32_t size, xm_u32_t *flags);
extern __stdcall xm_s32_t XM_seek_object(xmObjDesc_t objDesc, xm_u32_t offset, xm_u32_t whence);
extern __stdcall xm_s32_t XM_ctrl_object(xmObjDesc_t objDesc, xm_u32_t cmd, void *arg);
```

Listing 5.32: /user/libxm/include/xmhypercalls.h

The following services are implemented through the object interface:

- Communication ports.
- Console output.
- Health monitoring logs.
- Memory access.
- XtratuM and partition status.
- Trace logs.
- Serial ports.

1170

1175

For example, the XM_hm_status() hypercall is implemented in the libxm as:

```
xm_s32_t XM_hm_status(xmHmStatus_t *hmStatusPtr) {
    if (!(libXmParams.partCtrlTab->flags&XM_PART_SYSTEM))
        return XM_PERM_ERROR;
    if (!hmStatusPtr) {
        return XM_INVALID_PARAM;
    }
    return XM_ctrl_object(OBJDESC_BUILD(OBJ_CLASS_HM, XM_HYPERVISOR_ID, 0), XM_HM_GET_STATUS,
        hmStatusPtr);
}
```

Listing 5.33: /user/libxm/common/hm.c

5.20 Manpages summary

Below is a summary of the manpages. A detailed information is provided in the document “Volume 4: Reference Manual”.

Hypercall	Description
<code>XM_clear_irqmask</code>	Unmask interrupts.
<code>XM_clear_irqpend</code>	Unmask interrupts.
<code>XM_create_queuing_port</code>	Create a queuing port.
<code>XM_create_sampling_port</code>	Create a sampling port.
<code>XM_disable_irqs</code>	Replaced by <code>XM_sparcv8_clear_pil()</code> on Sparc processors.
<code>XM_enable_irqs</code>	Replaced by <code>XM_sparcv8_set_pil()</code> on Sparc processors.
<code>XM_get_partition_mmap</code>	Return a pointer to the memory map table (MMT).
<code>XM_get_partition_status</code>	Get the current status of a partition.
<code>XM_get_plan_status</code>	Return information about the scheduling plans.
<code>XM_get_queuing_port_info</code>	Get the info of a queuing port.
<code>XM_get_queuing_port_status</code>	Get the status of a queuing port.
<code>XM_get_sampling_port_info</code>	Get the info of a sampling port.
<code>XM_get_sampling_port_status</code>	Get the status of a sampling port.
<code>XM_get_system_status</code>	Get the current status of the system.
<code>XM_get_time</code>	Retrieve the time of the specified clock.
<code>XM_halt_partition</code>	Terminates a partition.
<code>XM_halt_system</code>	Stop the system.
<code>XM_hm_open</code>	Open the health monitoring log stream.
<code>XM_hm_read</code>	Read a health monitoring log entry.
<code>XM_hm_seek</code>	Sets the read position in the health monitoring stream.
<code>XM_hm_status</code>	Get the status of the health monitoring log stream.
<code>XM_idle_self</code>	Idles the execution of the calling partition.
<code>XM_mask_irq</code>	Obsoleted by <code>XM_set_irqmask()</code> .
<code>XM_memory_copy</code>	Copy copies data from/to address spaces.
<code>XM_multicall</code>	Execute a sequence of hypercalls.
<code>XM_params_get_PCT</code>	Return the address of the PCT.
<code>XM_read_console</code>	Print a string in the hypervisor console.
<code>XM_read_sampling_message</code>	Reads a message from the specified sampling port.
<code>XM_receive_queuing_message</code>	Receive a message from the specified queuing port.
<code>XM_request_irq</code>	Request to receive an interrupt.
<code>XM_reset_partition</code>	Reset a partition.
<code>XM_reset_system</code>	Reset the system.
<code>XM_resume_partition</code>	Resume the execution of a partition.
<code>XM_route_irq</code>	Link an interrupt with the vector which is generated when the interrupt is issued.
<code>XM_send_queuing_message</code>	Send a message in the specified queuing port.
<code>XM_set_irqmask</code>	Mask interrupts.
<code>XM_set_irqpend</code>	Mask interrupts.
<code>XM_set_page_type</code>	Changes the type of the physical page <code>pAddr</code> to <code>type</code> .
<code>XM_set_plan</code>	Request a plan switch at the end of the current MAF.
<code>XM_set_timer</code>	Arm a timer.
<code>XM_shutdown_partition</code>	Send a shutdown interrupt to a partition.
<code>XM_sparcv8_atomic_add</code>	Atomic add.
<code>XM_sparcv8_atomic_and</code>	Atomic bitwise AND.
<code>XM_sparcv8_atomic_or</code>	Atomic bitwise OR.
<code>XM_sparcv8_clear_pil</code>	Clear the PIL field of the PSR (allows interrupts).
<code>XM_sparcv8_flush_cache</code>	Flush data cache. Assembly hypercall.
<code>XM_sparcv8_flush_regwin</code>	Save the contents of the register window.
<code>XM_sparcv8_flush_tlb</code>	The TLB cache is invalidated. Assembly hypercall.
<code>XM_sparcv8_get_flags</code>	Replaced by <code>XM_sparcv8_get_psr()</code> on Sparc processors.
<code>XM_sparcv8_get_psr</code>	Get the ICC and PIL flags from the virtual PSR processor register.
<code>XM_sparcv8_inport</code>	Read from a hardware I/O port.
<code>XM_sparcv8_iret</code>	Return from an interrupt.
<code>XM_sparcv8_outport</code>	Write in a hardware I/O port.

Hypercall	Description
<code>XM_sparcv8_set_flags</code>	Set the ICC flags on the PSR processor register.
<code>XM_sparcv8_set_pil</code>	Set the PIL field of the PSR (disallow interrupts).
<code>XM_sparcv8_set_psr</code>	Set the ICC and PIL flags on the virtual PSR processor register.
<code>XM_suspend_partition</code>	Suspend the execution of a partition.
<code>XM_trace_event</code>	Records a trace entry.
<code>XM_trace_open</code>	Open a trace stream.
<code>XM_trace_read</code>	Read a trace event.
<code>XM_trace_seek</code>	Sets the read position in a trace stream.
<code>XM_trace_status</code>	Get the status of a trace stream.
<code>XM_unmask_irq</code>	Obsoleted by <code>XM_clear_irqmask()</code> .
<code>XM_update_page32</code>	Writes <code>val</code> in <code>pAddr</code> .
<code>XM_write_console</code>	Print a string in the hypervisor console.
<code>XM_write_register32</code>	Modify a processor control register.
<code>XM_write_sampling_message</code>	Writes a message in the specified sampling port.

This page is intentionally left blank.

Chapter 6

Binary Interfaces

This section covers the data types and the format of the files and data structures used by XtratuM.

Only the first section, describing the data types, is needed for the partition developer. The remaining sections contain material for advanced users. The `libxm.a` library provides a friendly interface that hides most of the low level details explained in this chapter.

1180

6.1 Data representation

The data types used in the XtratuM interfaces are compiler and machine cross development independent. This is specially important when manipulating the configuration files. These files may be created in a little-endian system (like the PC) while LEON2 is a big-endian one.



1185

XtratuM follows the next conventions:

Unsigned	Signed	Size (bytes)	Alignment (bytes)
<code>xm_u8_t</code>	<code>xm_s8_t</code>	1	1
<code>xm_u16_t</code>	<code>xm_s16_t</code>	2	4
<code>xm_u32_t</code>	<code>xm_s32_t</code>	4	4
<code>xm_u64_t</code>	<code>xm_s64_t</code>	8	8

Table 6.1: Data types.

These data types has to be stored in big-endian order, that is, the most significant byte standing at the lower address (0x..00) and the least significant byte standing to the upper address (0x..03).

The “C” declaration which meets these definitions is presented in the next listing:

```
// Basic types
typedef unsigned char xm_u8_t;
typedef char xm_s8_t;
typedef unsigned short xm_u16_t;
typedef short xm_s16_t;
typedef unsigned int xm_u32_t;
typedef int xm_s32_t;
typedef unsigned long long xm_u64_t;
typedef long long xm_s64_t;
```

Listing 6.1: `/core/include/arch/arch.types.h`

For future compatibility, most data structures contain version information. It is a `xm_u32_t` data type with 3 fields: version, subversion and revision. The following macros can be used to manipulate those fields:

```
#define XM_SET_VERSION(_ver, _subver, _rev) ((((_ver)&0xFF)<<16)|(((
    _subver)&0xFF)<<8)|(((_rev)&0xFF))
#define XM_GET_VERSION(_v) (((_v)>>16)&0xFF)
#define XM_GET_SUBVERSION(_v) (((_v)>>8)&0xFF)
#define XM_GET_REVISION(_v) ((_v)&0xFF)
```

Listing 6.2: `/core/include/xmef.h`

6.2 Hypercall mechanism

1190 An hypercall is implemented by a trap processor instruction that transfers the control to XtratuM code, and sets the processor in supervisor mode.

There are two kind of hypercalls: normal and assembly. Each type of hypercall use a different trap number:

```
#define XM_HYPERCALL_TRAP 0xF0
#define XM_ASMHYPERCALL_TRAP 0xF1
```

Listing 6.3: `/core/include/arch/xm_def.h`

1195 The `XM_ASMHYPERCALL_TRAP` hypercall entry is needed for the `XM_sparcv8_flush_regwin()`, `XM_sparcv8_iret()` and `XM_sparcv8_get_flags()` calls. In this case, the XtratuM entry code does not prepare the processor to execute “C” code.

6.3 Executable formats overview

XtratuM core does not have the capability to “load” partitions. It is assumed that when XtratuM starts its execution, all the partition code and data required to execute each partition is already in main memory. Therefore, XtratuM does not contain code to manage executable images. The only information required by XtratuM to execute a partition is the address of the partition image header (`xmImageHdr`).

1200 The partition images, as well as the XtratuM image, shall be loaded by a resident software, which acts as the boot loader.

The *XEF* (XtratuM Executable Format) has been designed as a robust format to copy the partition code (and data) from the partition developer to the final target system.

1205 The XtratuM image shall also be in XEF format. From the resident software point of view, XtratuM is just another image that has to be copied into the appropriate memory area.

The main features of the XEF format are:

- Simpler than the ELF. The ELF format is a rich and powerful specification, but most of its features are not required.
- Content checksum. Which allows to detect transmission errors.
- 1210 • Compress the content. This feature greatly reduce the space of the image; consequently the deploy time.

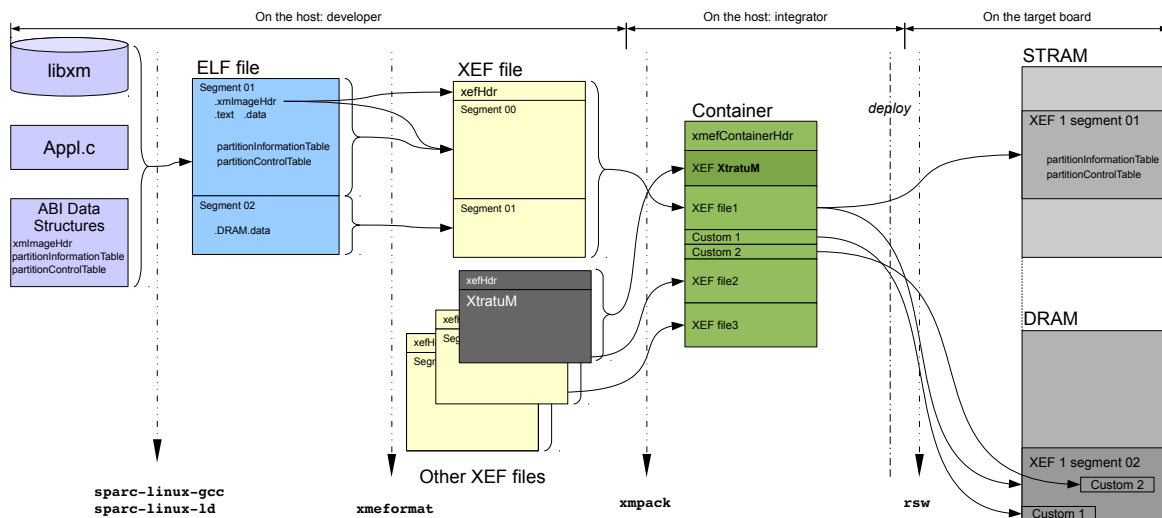


Figure 6.1: Executable formats.

- Encrypt the content. Not implemented.
- Partitions can be placed in several non-contiguous memory areas.

The *container* is a file which contains a set of XEF files. It is like a tar file (with important internal differences). The resident software shall be able to manage the container format to extract the partitions (XEF files); and also the XEF format to copy them to the target memory addresses. 1215

The signature fields, are constants used to identify and locate the data structures. The value that shall contain these fields on each data structure is defined right above the corresponding declaration.

6.4 Partition ELF format

A *partition image* contains all the information needed to “execute” the partition. It does not have loading or booting information. It contains one *image header structure*, one or more *partition header structures*, as well as the code and data that will be executed. 1220

Since multiple partition headers is an experimental feature (to support multiprocessor in a partition), we will assume in what follows that a partition file contains only one image header structure and one partition header structure.

Note: all the addresses of partition image are absolute addresses which refer to the target RAM memory locations. 1225

6.4.1 Partition image header

The partition image header is a data structure with the following fields:

```
struct xmImageHdr {
#define XMEF_PARTITION_MAGIC 0x24584d69 // $XMi
    xm_u32_t sSignature;
    xm_u32_t compilationXmAbiVersion; // XM's abi version
    xm_u32_t compilationXmApiVersion; // XM's api version
}
```

```

    xm_u32_t noCustomFiles;
    struct xefCustomFile customFileTab[CONFIG_MAX_NO_CUSTOMFILES];
    xm_u32_t eSignature;
};

```

Listing 6.4: /core/include/xmef.h

sSignature and eSignature: Holds the start and end signatures which identifies the structure as a XtratuM partition image.

1230 compilationXmAbiVersion: XtratuM ABI version used to compile the partition. That is, the ABI version of the libxm and other accompanying utilities used to build the XEF file.

compilationXmApiVersion: XtratuM API version used to compile the partition. That is, the API version of the libxm and other accompanying utilities used to build the XEF file.

The current values of these fields are:

```

#define XM_ABI_VERSION 3
#define XM_ABI_SUBVERSION 1
#define XM_ABI_REVISION 0

#define XM_API_VERSION 3
#define XM_API_SUBVERSION 1
#define XM_API_REVISION 2

```

Listing 6.5: /core/include/hypercalls.h

Note that these values may be different to the API and ABI versions of the running XtratuM. This information is used by XtratuM to check that the partition image is compatible.

1235 noCustomFiles: The number of extra files accompanying the image. If the image were Linux, then one of the modules would be the *initrd* image. Up to CONFIG_MAX_NO_FILES can be attached. The moduleTab table contains the locations in the RAM's address space of the partition where the modules shall be copied (if any). See section 5.18.

customFileTab: Table information about the customisation files.

```

struct xefCustomFile {
    xmAddress_t sAddr;
    xmSize_t size;
};

```

Listing 6.6: /core/include/xmef.h

sAddr: Address where the customisation file shall be loaded.

1240 size: Size of the customisation file.

The address where the custom files are loaded shall belong to the partition.



The xmImageHdr structure has to be placed in a section named “.xmImageHdr”. An example of how the header of a partition can be created is shown in section 5.4.

1245 The remainder of the image is free to the partition developer. There is not a predefined format or structure of where the code and data sections shall be placed.

6.4.2 Partition control table (PCT)

In order to minimize the overhead of the para-virtualised services, XtratuM defines a special data structure which is shared between the hypervisor and the partition called *Partition control table* (PCT). There is a PCT for each partition. XtratuM uses the PCT to send relevant operating information to the partitions. The partition is only allowed to read.

```
typedef struct {
    xm_u32_t magic;
    xm_u32_t xmVersion; // XM version
    xm_u32_t xmAbiVersion; // XM's abi version
    xm_u32_t xmApiVersion; // XM's api version
    xm_u32_t resetCounter;
    xm_u32_t resetStatus;
    xm_u32_t cpuKhz;
    xmId_t id;
    xm_u32_t flags;
    xm_u32_t hwIrqs; // Hw interrupts belonging to the partition
    xm_s32_t noPhysicalMemAreas;
    xm_u8_t name[CONFIG_ID_STRING_LENGTH];
    xm_u32_t iFlags;
    xm_u32_t hwIrqsPend; // pending hw irqs
    xm_u32_t hwIrqsMask; // masked hw irqs

    xm_u32_t extIrqsPend; // pending extended irqs
    xm_u32_t extIrqsMask; // masked extended irqs

    struct pctArch arch;
    struct {
        xm_u32_t noSlot:16, reserved:16;
        xm_u32_t id;
        xm_u32_t slotDuration;
    } schedInfo;
    xm_u16_t trap2Vector[NO_TRAPS];
    xm_u16_t hwIrq2Vector[CONFIG_NO_HWIRQS];
    xm_u16_t extIrq2Vector[XM_VT_EXT_MAX];
} partitionControlTable_t;
```

Listing 6.7: /core/include/guest.h

The libxm call `XM_params_get_PCT()` returns a pointer to the PCT.

1250

The architecture dependent part is defined in:

```
struct pctArch {
    xmAddress_t tbr;
#ifdef CONFIG_MMU
    xmAddress_t ptdL1;
    xm_u32_t faultStatusReg;
    xm_u32_t faultAddressReg;
#endif
};
```

Listing 6.8: /core/include/arch/guest.h

signature: Signature to identity this data structure as a PIT.

xmAbiVersion: The Abi version of the currently running XtratuM. This value is filled by the running XtratuM.

1255 **xmApiVersion:** The Api version of the currently running XtratuM. This value is filled by the running XtratuM.

resetCounter: A counter of the number of partition resets. This counter is incremented when the partition is WARM reset. On a COLD reset it is set to zero.

resetStatus: If the partition had been reset by a `XM_reset_partition()` hypercall, then the value of the parameter `status` is copied in this field. Zero otherwise.

1260 **id:** The identifier of the partition. It is the unique number, specified in the `XM.CF` file, to unequivocally identify a partition.

hwIrqs: A bitmap of the hardware interrupts allocated to the partition. Hardware interrupts are allocated to the partition in the `XM.CF` file.

1265 **noPhysicalMemoryAreas:** The number of memory areas allocated to the partition. This value defines the size of the `physicalMemoryAreas` array.

name: Name of the partition.

hwIrqsPend: Bitmap of the hardware interrupts allocated to the partition delivered to the partition.

extIrqsPend: Bitmap of the extended interrupts allocated to the partition delivered to the partition.

hwIrqsMask: Bitmap of the extended interrupts allocated to the partition delivered to the partition.

1270 **extIrqsMask:**

In the current version there is no specific architecture data.

6.5 XEF format

The XEF is a wrapper for the files that may be deployed in the target system. There are three kind of files:

- Partition images.
- 1275 • The XtratuM image.
- Customisation files.

An XEF file has an header (see listing 6.9) and a set of *segments*. The segments, like in ELF, represent blocks of memory that will be loaded in RAM.

The tool `xmeformat` converts from ELF or plain data files to XEF format, see chapter 9.

```

struct xefHdr {
#define XEF_SIGNATURE 0x24584546
    xm_u32_t    signature;
    xm_u32_t    version;
#define XEF_DIGEST 0x1
#define XEF_COMPRESSED 0x4
#define XEF_RELOCATABLE 0x10
#define XEF_CONTENT_MASK 0xc0

```

```

#define XEF_CONTENT_HYPERVISOR 0x00
#define XEF_CONTENT_PARTITION 0x40
#define XEF_CONTENT_CUSTOMFILE 0x80
//#define XEF_CONTENT_RESERVED 0xc0

    xm_u32_t      flags;
    xm_u8_t       digest[XM_DIGEST_BYTES];
    xm_u8_t       payLoad[XM_PAYLOAD_BYTES];
    xm_u32_t      fileSize;
    xmAddress_t   segmentTabOffset;
    xm_s32_t      noSegments;
    xmAddress_t   customFileTabOffset;
    xm_s32_t      noCustomFiles;
    xmAddress_t   imageOffset;
    xm_u32_t      imageLength;
    xm_u32_t      deflatedImageLength;

    xmAddress_t   xmImageHdr;
    xmAddress_t   entryPoint;
#ifdef CONFIG_IA32
    xmAddress_t   relOffset;
    xm_s32_t      noRel;

    xmAddress_t   relaOffset;
    xm_s32_t      noRela;
#endif
};

```

Listing 6.9: /core/include/xmef.h

signature: A 4 bytes word to identify the file as an XEF format.

1280

version: Version of the XEF format.

flags: Bitmap of features present in the XEF image. It is a 4 bytes word. The existing flags are:

XEF_DIGEST: If set, then the `digest` field is valid and shall be used to check the integrity of the XEF file.

XEF_COMPRESSED: If set, then the partition binary image is compressed.

1285

XEF_CIPHERED: (**future extension**) to inform whether the partition binary is encrypted or not.

XEF_CONTENT: Specifies what kind of file is.

digest: when the `XEF_DIGEST` flag is set, this field holds the result of processing all the XEF file (supposing the `digest` field set to 0). The MD5 algorithm is used to calculate this field.

Despite the well known security flaws, we selected the MD5 digest algorithm because it has a reasonable trade-off between calculation time and the security level¹. Note that the `digest` field is used to detect not deliberate modifications rather than intentional attacks. In this scenario, the MD5 is a good choice.

1290

¹According to our tests, the time spent by more sophisticated digest algorithms such as SHA-2, Tiger or Whirlpool in the LEON3 processor was not acceptable. As illustration, 100 Kbytes took several seconds to be digested by a SHA-2 algorithm.

payload: This field holds 16 bytes which can freely be used by the partition supplier. It could be used to hold information such as partition's version, etc.

1295

The content of this field is used neither by XtratuM nor the resident software.

fileSize: XEF file size in bytes.

segmentTabOffset: Offset to the section table.

noSegments: Number of segments held in the XEF file. In the case of a customisation file, there will be only one segment.

1300

customFileTabOffset: Offset to the custom files table.

noCustomFiles: Number of custom files.

imageOffset: Offset to the partition binary image.

imageLength: Size of the partition binary image.

deflatedImageLength: When the XEF_COMPRESS flag is set, this field holds the size of the uncompressed partition binary image.

1305

xmImageHdr: Pointer to the partition image header structure (`xmImageHdr`). The `xmefformat` tool copies the address of the corresponding section in this file.

entryPoint: Address of the starting function.

Additionally, analogically to the ELF format, XEF contemplates the concept of *segment*, which is, a portion of code/data with a size and a specific load address. A XEF file includes a segment table (see listing 6.10) which describes each one of the sections of the image (custom data XEF files have only one section).

1310

```

struct xefSegment {
    xmAddress_t  startAddr;
    xm_u32_t     fileSize;
    xm_u32_t     deflatedFileSize;
    xmAddress_t  offset;
};

```

Listing 6.10: `/core/include/xmef.h`

startAddr: Address where the segment shall be located while it is being executed. This address is the one used by the linker to locate the image. If there is not MMU, then `physAddress=virtAddr`.

1315

fileSize: The size of the segment within the file. This size could be different from the memory required to be executed (for example a BSS usually requires more memory once loaded into memory).

deflatedFileSize: When the XEF_COMPRESS flag is set, this field holds the size of the segment when uncompressed.

1320

offset: Location of the segment expressed as an offset in the partition binary image.

6.5.1 Compression algorithm

The compression algorithm implemented is Lempel-Ziv-Storer-Szymanski (LZSS). It is a derivative of LZ77, that was created in 1982 by James Storer and Thomas Szymanski. A detailed description of the algorithm appeared in the article “Data compression via textual substitution” published in Journal of the ACM. 1325

The main features of the LZSS are:

1. Fairly acceptable trade-off between compression rate and decompression speed.
2. Implementation simplicity.
3. Patent-free technology.

Aside from LZSS, other algorithms which were regarded were: huffman coding, gzip, bzip2, LZ77, RLE and several combinations of them. Table 6.2 sketches the results of compressing XtratuM’s core binary with some of these compression algorithms. 1330

Algorithm	Compressed size	Compression rate (%)
LZ77	43754	44.20%
LZSS	36880	53.01%
Huffman	59808	23.80%
Rice 32bits	78421	0.10%
RLE	74859	4.60%
Shannon-Fano	60358	23.10%
LZ77/Huffman	36296	53.76%

Table 6.2: Outcomes of compressing the `xm_core.bin` (78480 bytes) file.

6.6 Container format

A *container* is a file which contains a set of XEF files.

The tool `xmpack` manages container files, see chapter 9.

A *component* is an executable binary (hypervisor or partition) jointly with associated data (configuration or customization file). The XtratuM component contains the files: `xm_core.bin` and `XM_CT.bin`. A partition component is formed by the partition binary file and zero or more customization files. 1335

XtratuM is not a boot loader. There shall be an external utility (the resident software or boot loader) which is in charge of coping the code and data of XtratuM and the partition from a permanent memory into the RAM. Therefore, the the container file is not managed by XtratuM but by the resident software, see chapter 7. 1340

Note also, that the container does not have information regarding where the components shall be loaded into RAM memory. This information is contained in the header of the binary image of each component.

The container file is like a packed filesystem which contains several the file metadata (name of the files) and the content of each file. Also, which file contains the executable image and the customisation data of each partition is specified. 1345

The container has the following elements:

- 1350 1. The header (`xmefContainerHdr` structure). A data structure which holds pointers (in the form of offsets) and the sizes to the remainder sections of the file.
2. The component table section, which contains an array of `xmefComponent` structures. Each element contains information of one component.
3. The file table section, which contains an array of files (`xmefFile` structure) in the container.
- 1355 4. The string table section. Contains the names of the files of the original executable objects. This is currently used for debugging.
5. The file data table section, with the actual data of the executable (XtratuM and partition images) and configuration files.

The container header has the following fields:

```

struct xmefContainerHdr {
    xm_u32_t    signature;
#define XM_PACKAGE_SIGNATURE 0x24584354 // $XCT
    xm_u32_t    version;
#define XMPACK_VERSION 3
#define XMPACK_SUBVERSION 0
#define XMPACK_REVISION 0
    xm_u32_t    flags;
#define XMEF_CONTAINER_DIGEST 0x1
    xm_u8_t     digest[XM_DIGEST_BYTES];
    xm_u32_t    fileSize;
    xmAddress_t partitionTabOffset;
    xm_s32_t    noPartitions;
    xmAddress_t fileTabOffset;
    xm_s32_t    noFiles;
    xmAddress_t strTabOffset;
    xm_s32_t    strLen;
    xmAddress_t fileDataOffset;
    xmSize_t    fileDataLen;
};

```

Listing 6.11: `/core/include/xmef.h`

signature: Signature field.

version: Version of the package format.

1360 **flags:**

digest: Not used. Currently the value is zero.

fileSize: The size of the container.

partitionTabOffset: The offset (relative to the start of the file) to the partition array section.

noPartitions: Number of partitions plus one (XtratuM is also a component) in the container.

1365 **componentOffset:** The offset (relative to the start of the file) to the component's array section.

fileTabOffset: The offset (relative to the start of the container file) to the files's array section.

noFiles: Number of files (XtratuM core, the XM_CT file, partition binaries, and partition-customization files) in the container.

strTabOffset The offset (relative to the start of the container file) to the strings table.

strLen The length of the strings table. This section contains all names of the files.

1370

fileDataOffset The offset (relative to the start of the container file) to the file data section.

fileDataLen The length of the file data section. This section contains all the contents of all the components.

Each entry of the partition table section describes all the XEF files that are part of each partition. Which contains the following fields:

```

struct xmefPartition {
    xm_s32_t    id;
    xm_s32_t    file;
    xm_u32_t    noCustomFiles;
    xm_s32_t    customFileTab[CONFIG_MAX_NO_CUSTOMFILES];
};

```

Listing 6.12: /core/include/xmef.h

id: The identifier of the partition.

file: The index into the file table section of the XEF partition image.

1375

noCustomFiles: Number of customisation files of this component, including.

customFileTab: List of custom file indexes.

The metadata of each file is store in the file table section:

```

struct xmefFile {
    xmAddress_t offset;
    xmSize_t size;
    xmAddress_t nameOffset;
};

```

Listing 6.13: /core/include/xmef.h

offset: The offset (relative to the start of the file data table section) to the data of this file in the container.

size: The size reserved to store this file. It is possible to define the size reserved in the container to store a file independently of the actual size of the file. See the section 9.3.1 tool.

1380

nameOffset: Offset, relative to the start of the strings table, of the name of the file.

The strings table contains the list of all the file names.

The file data section contains the data (with padding if `fileSize<=size`) of the files.

This page is intentionally left blank.

Chapter 7

Booting

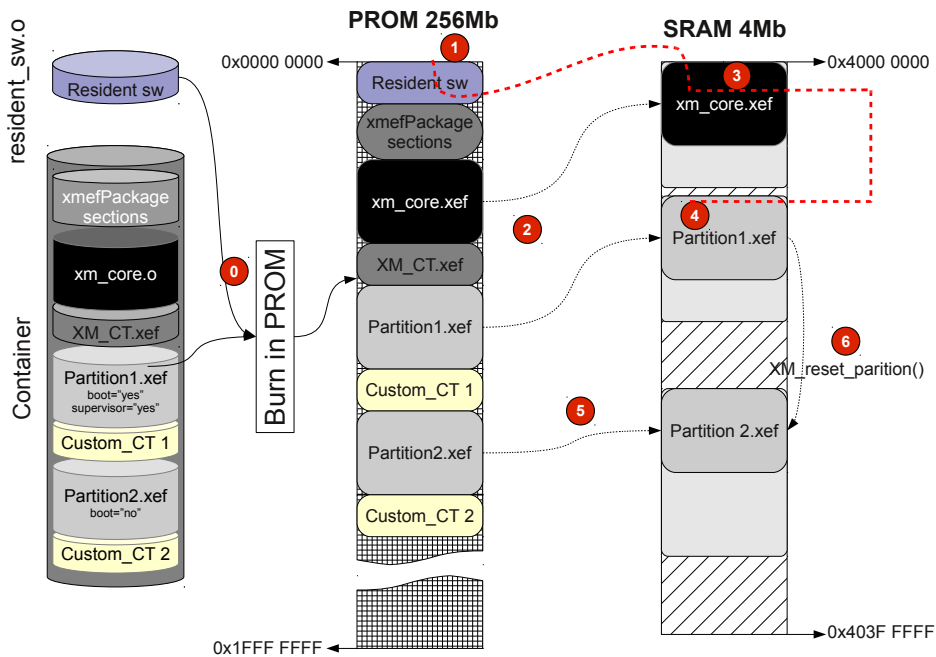


Figure 7.1: Booting sequence.

In the standard boot procedure of the LEON2 processor, the program counter register is initialized with the address 0x00000000. Contrarily to other computers, the PROM of the board does not have any kind of resident software-like booter¹ that takes the control of the processor after the reset. 1385

We have developed a small booting code called *resident software*, which is in charge of the initial steps of booting the computer. This software is not part of the container produced by the `xmpack` tool. It is prepended to the container by the burning script. 1390

The board has two kind of RAM memory: SRAM (4Mb) and SDRAM (128Mb).

¹Known as BIOS in the personal computer area.

7.1 Boot configuration

The *resident software* is in charge of loading into memory XtratuM, its configuration file (XM_CT) and any partition jointly with its customisation file as found in the container. The information hold by the XM_CT file is used to load any partition image. Additionally, the *resident software* informs XtratuM which partitions have to be booted.



After starting, XtratuM assumes that the partitions informed as ready-to-be-booted are in RAM/SRAM memory, setting them in running state right after it finishes its booting sequence

If a partition's code is not located within the container, then XtratuM sets the partition in HALT state until a system partition resets it by using `XM_reset_partition()` hypercall. In this case, the RAM image of the partition shall be loaded by a system partition through the `XM_memory_copy()` hypercall.

Note that there may be several booting partitions. All those partitions will be started automatically at boot time.

The boot sequence is sketched in figure 7.1.

① The deployment tool to burn the PROM of the board writes first the resident software and right after it the container, which should contain the XtratuM core and the booting partitions components. Note that the container can also contain the non-booting partitions.

① When the processor is started (reset) the resident software is executed. It is a small code that performs the following actions:

1. Initializes a stack (required to execute "C" code).
2. Installs a trap handler table (only for the case that its code would generate a fault, XtratuM installs a new trap handler during its initialisation).
3. Checks that the data following the resident software in the PROM is a container (a valid signature), and seeks the XtratuM hypervisor through container (a valid signature).
4. Copies the XtratuM hypervisor and booting partitions into RAM memory (②).
5. The address of the container (which contains the `ctCompTab`) is copied in the `%g2` processor register. Jumps to the entry point of the hypervisor in RAM memory.

③ XtratuM assumes no initial processor state. So, the first code has to be written in assembly (code/kernel/sparcv8/head.s), and performs the next actions: the `%g2` register is saved in a global "C" variable; general purpose registers are cleared; memory access rights are cleared; PSR², WIM³, TBR⁴ and Y⁵ processor control registers are initialized; sets up a working stack; and jumps to "C" code (code/kernel/setup.c).

The `setup()` function carries out the following actions:

1. Initializes the internal console.
2. Initializes the interrupt controller.
3. Detects the processor frequency (information extracted from the XML configuration file).
4. Initializes memory manager (enabling XtratuM to keep track of the use of the physical memory).

²PSR: Processor Status Register.

³WIM: Window Invalid Mask.

⁴TBF: Trap Base Register.

⁵Y: Extended data register for some arithmetic operations.

5. Initializes hardware and virtual timers.
6. Initializes the scheduler. 1430
7. Initializes the communication channels.
8. Booting partitions are set in NORMAL state and non-booting ones are set in HALT state.
9. Finally, the `setup` function calls the scheduler and becomes into the idle task.

- ④ Partition code is executed according to the configured plan.
- ⑤ A system partition can load from PROM or other source (serial line, etc.) the image of other partitions. 1435
- ⑥ The new ready partition is activated via a `XM_reset_partition()` service.

This page is intentionally left blank.

Chapter 8

Configuration

This section describes how XtratuM is configured. There are two levels of configuration. A first level which affects the source code to customise the resulting XtratuM executable image. Since XtratuM does not use dynamic memory to setup internal data structures, most of these configuration parameters are related to the size, or ranges, of the statically created data structures (maximum number of partitions, channels, etc..).

1440

The second level of configuration is done via an XML file. This file configures the resources allocated to each partition.

8.1 XtratuM source code configuration (menuconfig)

The first step in the XtratuM configuration is to configure the source code. This task is done using the same tool than the one used in Linux, which are commonly called “make menuconfig”.

1445

There are two different blocks that shall be configured: 1) XtratuM source code; and 2) the resident software. The configuration menu of each block is presented one after the other when executed the “\$ make menuconfig” from the root source directory. The selected configuration are stored in the files `core/.config` and `/user/bootloaders/rsw/.config` for XtratuM and the resident software respectively.

1450

The next table lists all the XtratuM configuration options and its default values. Note that since there are logical dependencies between some options, the menuconfig tool may not show all the options. Only the options that can be selected are presented to the user.

There has been major changes from version 2.2 to 3.x in the source code configuration.

1455



Parameter	Type	Default value
Processor		
SPARC cpu	choice	[Leon3]
Board	choice	[TSim] [GR-CPCI-XC4V] [GR-XC3S-1500]
SPARC memory protection schema	choice	[MMU]
Support AMBA bus PnP	bool	y
Enable VMM update hypercalls	bool	y
Enable cache	bool	y
Enable cache snoop	bool	n
<i>Continues...</i>		

Parameter	Type	Default value
Enable instruction burst fetch	bool	n
Flush cache after context switch	bool	n
Physical memory layout		
XM load address	hex	0x40000000
Debug and profiling support	bool	y
Max. identifier length (B)	int	32
Hypervisor		
Kernel stack size (KB)	int	8
System integrity	bool	y
Context switch threshold (usec)	int	200
Drivers		
Enable UART rx interrupt	bool	n
Enable UART flow control	bool	n
DSU samples UART port	bool	n
Enable VGA support	bool	n
Objects		
Console initial buffer length	int	256
Enable XM/partition status accounting	bool	n

Sparc cpu: Processor model.

Board: Enables the specific board features: write protection units, timers, UART and interrupt controller.

SPARC memory protection schema: Select the processor mechanism that will be used to implement memory protection. If MMU is available then it is the best choice. With WPR, only write protection can be enforced.

Enable cache: If selected, the processor cache is enabled. All partitions will be executed with cache enabled unless explicitly disabled on each partition through the XM_CF file.

If this option is not selected, the cache memory is disabled. Neither XtratuM nor the partitions will be able to use the cache.

Flush cache after context switch: Forces a cache flush after a partition context switch.

DSU samples UART port: If the UART port used to print console messages is also used by the DSU (Debugging Support Unit), then this option shall be set.

If the DSU is used, then the control bits of the UART does not change. In this case, bytes are sent with a timeout.

XtratuM load address: Physical RAM address where XtratuM shall be copied. This value shall be the same than the one specified in the XM_CF file.

Console initial buffer length: Size of the internal console buffer. This buffer it used to store the messages written by XtratuM or by partitions using the XM_console_write() hypercall. The larger the buffer, the lower the chances of loosing data.

Enable voluntary preemption support:

Enable experimental features: Enable this option to be able to select experimental ones. This option does not do anything on itself, just shows the options marked as experimental.

Kernel stack size (KB): Size of the stack allocated to each partition. It is the stack used by XtratuM when attending the partition hypercalls. 1480

Do not change (reduce) this value unless you know what you are doing.

Debug and profiling support: XtratuM is compiled with debugging information (gcc flag “-ggdb”) and assert code is included. This option should be used only during the development of the XtratuM hypervisor.

Maximum identifier length (B): The maximum string length (including the terminating “0x0” character) of the names: partition name, port name, plan, etc. Since the names are only used for debugging, 16 characters is a fair number. 1485

GCoverage support: Experimental.

Enable UART support: If enabled, XtratuM will use the UART to output console messages; otherwise the UART can be used by a partition. 1490

Enable XM/partition status accounting: Enable this option to collect statistical information of XtratuM itself and partitions.

Note that this feature increases the overhead of most of the XtratuM operations.

8.2 Resident software source code configuration (menuconfig)

The resident software (RSW) configuration parameters are hard-coded in the source code in order to generate a self-contained stand alone executable code. 1495

After the configuration of the XtratuM source code, the “\$ make menuconfig” shows the RWS configuration menu. The selected configuration is stored in the file `user/bootloaders/rsw/.config`.

The following parameters can be configured:

Parameter	Type	Default value
RSW memory layout		
Read-only section addresses	hex	0x40200000
Read/write section addresses	hex	0x40090000
CPU frequency (KHz)	int	50000
Enable UART support	choice	[UART1] [UART2]
UART baud rate	int	115200
Stack size (KB)	int	8
Stand-alone version	bool	no
Container physical location address	hex	0x4000

Stack size (KB):

Read-only section addresses: RSW memory layout. Read-only area. The resident software will be compiled to use this area. 1500

Read/write section addresses: RSW memory layout. Read/write area used by the resident software.

CPU frequency (KHz): The **processor frequency** is passed to XtratuM via the `XM_CF` file, but in the case of the RSW it has to be specified in the source code, since it has no run-time configuration. The processor frequency is used to configure the UART clock. 1505

Enable UART support: Select the serial line to write messages to.

UART baud rate: The baud rate of the UART. Note that the baud rate used by XtratuM is configured in the XM.CF file, and not in the source code configuration.

Stand-alone version: If not set, then the resident software shall be linked jointly with the container. That is, the final resident software image shall contain, as data, the container.

If set, then the container is not linked with the resident software. The address where the container will be copied in RAM is specified by the next option:

Container physical location address: Address of the container. In case of the stand-alone version.

8.2.1 Memory requirements

The memory footprint of XtratuM is independent of the workload (number of partitions, channels, etc.) The memory needed depends only on actual workload defined in the XM.CF file. The size of the compiled configuration provides an accurate estimation of the memory what will use XtratuM to run it. Note that it is not the size of the file, but the memory required by all the sections (including those not allocatable ones: `.bss`).

The resident software can be executed in ROM or RAM memory (if the ROM technology allows to run eExecute-in-Place XiP code). The resident software has no initialised data segment, only the `.text` and `.rodata` segments are required (the `.got` and `.eh_frame` segments are not used). The memory footprint of the RSW depends on whether the debug messages are enabled or not; and it can be obtained from the final resident software code (the file created with the `build_rsw` helper utility) using the `readelf` utility.

The next example shows where the size of the RSW code is printed (column labeled **MemSiz**):

```
# sparc-linux-readelf -l resident_sw
Elf file type is EXEC (Executable file)
Entry point 0x4020102c
There are 5 program headers, starting at offset 52

Program Headers:
Type      Offset    VirtAddr  PhysAddr  FileSiz  MemSiz   Flg Align
LOAD      0x0000d4 0x00004000 0x00004000 0x196f0 0x196f0  RW 0x1 // 00 segment
LOAD      0x0197c4 0x40090000 0x0001d6f0 0x00048 0x00048  RW 0x4 // 01 segment
LOAD      0x019808 0x40090048 0x40090048 0x00000 0x02000  RW 0x8 // 02 segment
LOAD      0x019810 0x40200000 0x40200000 0x01f44 0x01f44  R E 0x8 // 03 segment
GNU_STACK 0x000000 0x00000000 0x00000000 0x00000 0x00000  RWE 0x4 // 04 segment

Section to Segment mapping:
Segment Sections...
00  .container
01  .got .eh_frame
02  .bss
03  .text .rodata
04
```

Listing 8.1: Resident software memory footprint example.

The next table summarises the size of the resident software for different configurations:

Board type	Debug messages	size
LEON3	disabled	0x01f44
LEON3	enabled	0x01ff4
LEON2	disabled	0x01f44
LEON2	enabled	0x01ff4

8.3 Hypervisor configuration file (XM_CF)

The XM_CF file defines the system resources, and how they are allocated to each partition.

For an exact specification of the syntax (mandatory/optional elements and attributed, and how many times an element can appear) the reader is referred to the XML schema definition in the Appendix A.

8.3.1 Data representation and XPath syntax

When representing physical units, the following syntax shall be used in the XML file:

1530

Time: Pattern: “[0-9]+(. [0-9]+)?([mu]?[sS])”

Examples of valid times:

```
9s      # nine seconds.
10ms    # ten milliseconds.
0.5ms   # zero point five milliseconds.
500us   # five hundred microseconds =0.5ms
```

1535

Size: Pattern: “[0-9]+(. [0-9]+)?([MK]?B)”

Examples of valid sizes:

```
90B     # ninety bytes.
50KB    # fifty Kilo bytes =(50*1024) bytes.
2MB     # two mega bytes =(2*1024*1024) bytes.
2.5KB   # two point five kilo bytes =2560B.
```

1540

It is advised not to use the decimal point on sizes.

Frequency: Pattern: “[0-9]+(. [0-9]+)?([MK] [Hh]z)”

Examples of valid frequencies:

```
80Mhz   # Eighty mega hertz = 80000000 hertz.
20000Khz # Twenty mega hertz = 20000000 hertz.
```

1545

Boolean: Valid values are: “yes”, “true”, “no”, “false”.

Hexadecimal: Pattern: “0x[0-9a-fA-F]+”

Examples of valid numbers:

```
0xFfffFfff, 0x0, 0xF1, 0x80
```

1550

An XML file is organised as a set of nested elements, each element may contain attributes. The XPath syntax is used to refer to the objects (elements and attributes). Examples:

/SystemDescription/PartitionTable The element PartitionTable contained inside the element SystemDescription, which is the root element (the starting slash symbol).

1555

/SystemDescription/@name Refers to the attribute **./@name** of the element **SystemDescription**.

./Trace/@bitmask Refers to the attribute **./@bitmask** of a **./Trace** element. The location of the element **./Trace** in the xml element hierarchy is relative to the context where the reference appears.

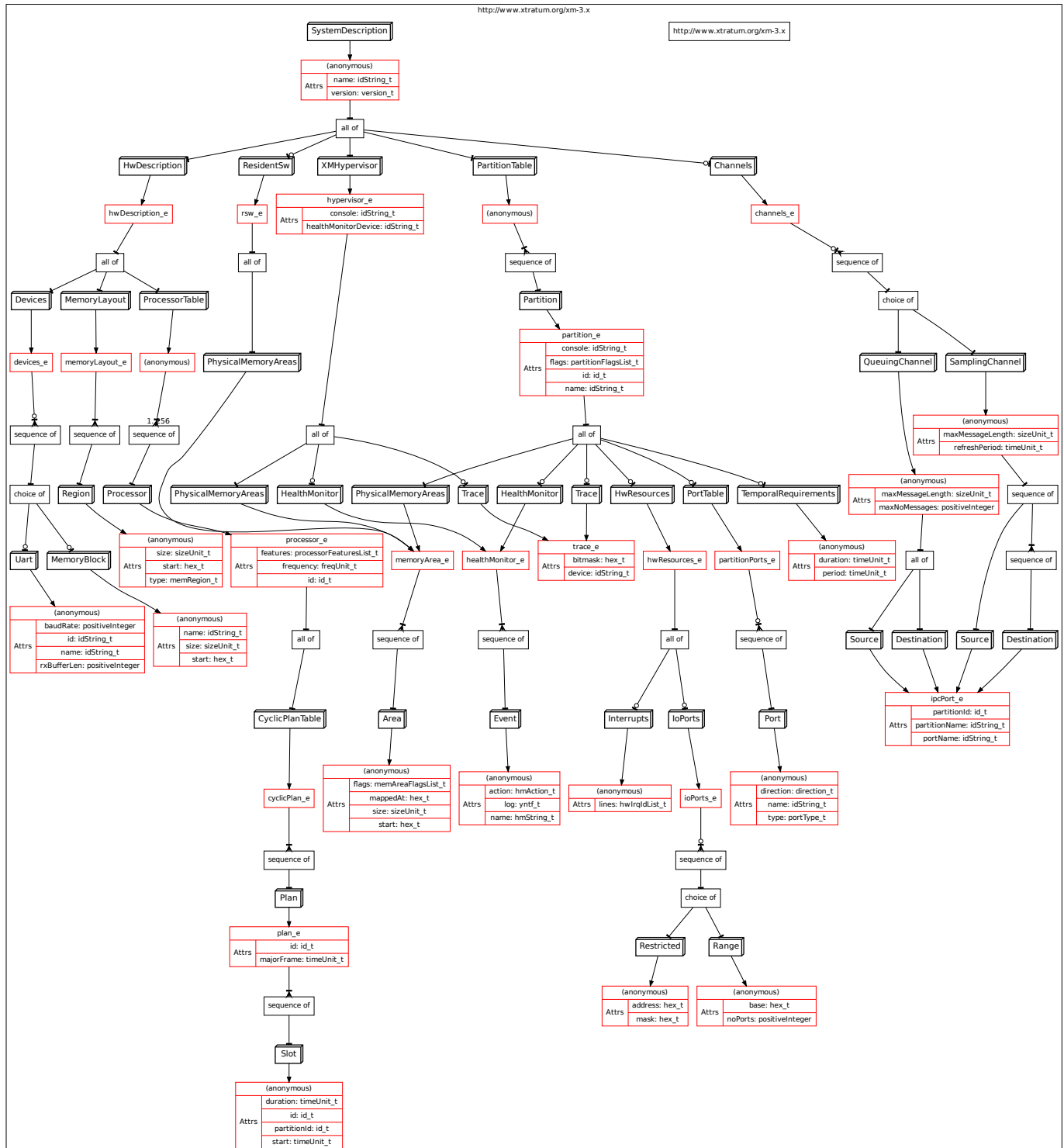


Figure 8.1: XML Schema diagram.

8.3.2 The root element: /SystemDescription

Figure 8.1 is a graphical representation of the schema of the XML configuration file. The types of the attributes are not represented, see the appendix A for the complete schema specification. An arrow ended with circle are optional elements. 1560

Figure 8.2 on page 86 is a compact graphical representation of the nested structure a sample XM_CF configuration file (the listing A.2 is the actual xml file for this representation). Solid-lined boxes represent elements. Dotted boxes contain attributes. The nested boxes represent the hierarchy of elements. 1565

The root element is `"/SystemDescription"`, which contain the mandatory `./@version`, `./@name` and `./@xmlns` attributes. The xmlns name space shall be `"http://www.xtratum.org/xm-2.3"`.

There are five second-level elements:

`/SystemDescription/XMHypervisor` Specifies the board resources (memory, and processor plan) and the hypervisor health monitoring table. 1570

`/SystemDescription/ResidentSw` This is an optional element which for providing information to XtratuM about the resident software.

`/SystemDescription/PartitionTable` This is a container element which holds all the `./partition` elements.

`/SystemDescription/Channels` A sequence of channels which define port connections. 1575

`/SystemDescription/HwDescription` Contain the configuration of physical and virtual resources.

8.3.3 The /SystemDescription/XMHypervisor element


There are two optional attributes `./@console` and `./@healthMonitoringDevice`. The values of these attributes shall be the name of a device defined in the `/SystemDescription/HwDescription/Devices` section.

Mandatory elements:

`./PhysicalMemoryAreas` Sequence of memory areas allocated to XtratuM. 1580

Optional elements:

`./HealthMonitoring` Contains a sequence of health monitoring event elements.

Not all HM actions can be associated with all HM events. Consult the allowed actions in the "Volume 4: Reference Manual". 1585 

`./Trace` Defines where to store the traces messages emitted by XtratuM (the value of the attribute `./@device` shall be a the name of a device defined in `/SystemDescription/Devices`); and the hexadecimal bit mask to filter out which traces will not be stored (`./@bitmask`).

A health monitoring event element contains the following attributes:

`./event/@name` The event's name. Below is the list of available events:

```
<xs:simpleType name="hmString_t">
  <xs:restriction base="xs:string">
    <xs:enumeration value="XM_HM_EV_INTERNAL_ERROR"/>
    <xs:enumeration value="XM_HM_EV_UNEXPECTED_TRAP"/>
    <xs:enumeration value="XM_HM_EV_PARTITION_ERROR"/>
    <xs:enumeration value="XM_HM_EV_PARTITION_INTEGRITY"/>
  </xs:restriction>
</xs:simpleType>
```

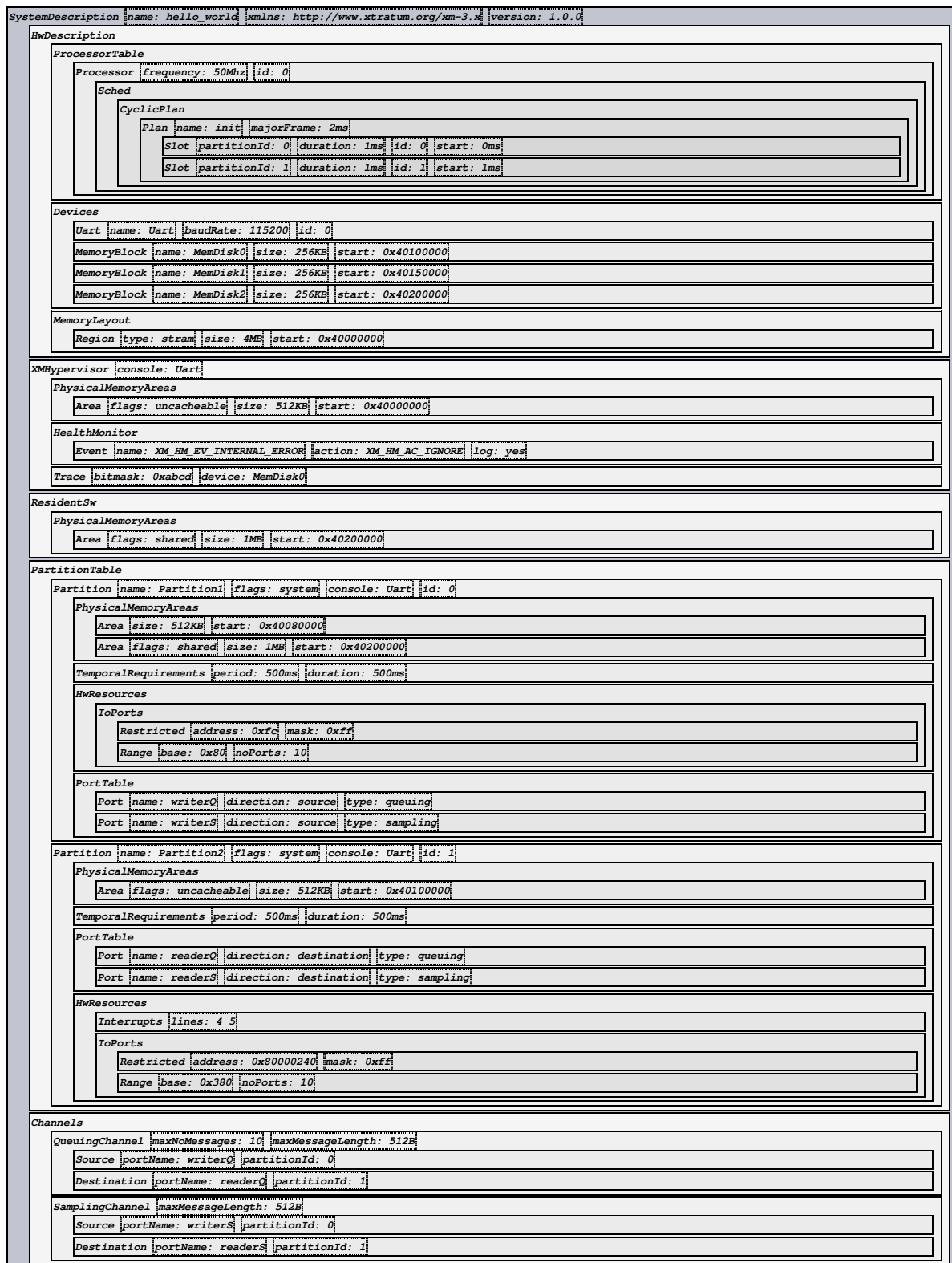


Figure 8.2: Graphical view of an example XM.CF configuration file (see the XML file in section A.2).


```

<xs:enumeration value="XM_HM_EV_MEM_PROTECTION"/>
<xs:enumeration value="XM_HM_EV_OVERRUN"/>
<xs:enumeration value="XM_HM_EV_SCHED_ERROR"/>
<xs:enumeration value="XM_HM_EV_WATCHDOG_TIMER"/>
<xs:enumeration value="XM_HM_EV_INCOMPATIBLE_INTERFACE"/>
<xs:enumeration value="XM_HM_EV_WRITE_ERROR"/>
<xs:enumeration value="XM_HM_EV_INSTR_ACCESS_MMU_MISS"/>
<xs:enumeration value="XM_HM_EV_INSTR_ACCESS_ERROR"/>
<xs:enumeration value="XM_HM_EV_UNIMPLEMENTED_FLUSH"/>
<xs:enumeration value="XM_HM_EV_WATCHPOINT_DETECTED"/>
<xs:enumeration value="XM_HM_EV_DATA_ACCESS_ERROR"/>
<xs:enumeration value="XM_HM_EV_DATA_ACCESS_MMU_MISS"/>
<xs:enumeration value="XM_HM_EV_INSTR_ACCESS_EXCEPTION"/>
<xs:enumeration value="XM_HM_EV_ILLEGAL_INSTR"/>
<xs:enumeration value="XM_HM_EV_PRIVILEGED_INSTR"/>
<xs:enumeration value="XM_HM_EV_FP_DISABLED"/>
<xs:enumeration value="XM_HM_EV_CP_DISABLED"/>
<xs:enumeration value="XM_HM_EV_REGISTER_HARDWARE_ERROR"/>
<xs:enumeration value="XM_HM_EV_MEM_ADDR_NOT_ALIGNED"/>
<xs:enumeration value="XM_HM_EV_FP_EXCEPTION"/>
<xs:enumeration value="XM_HM_EV_CP_EXCEPTION"/>
<xs:enumeration value="XM_HM_EV_DATA_ACCESS_EXCEPTION"/>
<xs:enumeration value="XM_HM_EV_TAG_OVERFLOW"/>
<xs:enumeration value="XM_HM_EV_DIVIDE_EXCEPTION"/>
</xs:restriction>
</xs:simpleType>

```

Listing 8.2: /user/tools/xmcparser/xmc.xsd.in

./event/@action The name of the action associated with this event. Below in the list of available actions:

```

<xs:simpleType name="hmAction_t">
  <xs:restriction base="xs:string">
    <xs:enumeration value="XM_HM_AC_IGNORE"/>
    <xs:enumeration value="XM_HM_AC_SHUTDOWN"/>
    <xs:enumeration value="XM_HM_AC_COLD_RESET"/>
    <xs:enumeration value="XM_HM_AC_WARM_RESET"/>
    <xs:enumeration value="XM_HM_AC_SUSPEND"/>
    <xs:enumeration value="XM_HM_AC_HALT"/>
    <xs:enumeration value="XM_HM_AC_PROPAGATE"/>
    <xs:enumeration value="XM_HM_AC_SWITCH_TO_MAINTENANCE" />
  </xs:restriction>
</xs:simpleType>

```

Listing 8.3: /user/tools/xmcparser/xmc.xsd.in

./event/@log Boolean flag to select whether the event will be logged or not.

1590

8.3.4 The /SystemDescription/HwDescription element

It contains three mandatory elements:

./HwDescription/ProcessorTable Which holds a sequence of **./Processor** elements. Each processor element describes one physical processor: the processor clock **./@frequency** (the frequency units has to be specified), **./@id** (zero in a mono-processor system), and an optional **./@features** attribute. The **./@features** attribute contains a list of specific processor features than can be selected. Currently, only the memory protection workaround (“XM_CPU_LEON2_WA1”), for the memory mapped processor registers bug¹.

1595

¹Some key processor registers, needed to guarantee the spatial isolation, are mapped memory addresses which are not monitored/protected by the write protection mechanism. The workaround consists in protecting this register area using the watchpoint mechanism. The workaround is only applicable if the watchpoint facility is present.

Also, the `./ProcessorTable/Processor` element defines the scheduling plan of this processor. It is specified in the element `./Processor/Sched/CyclicPlan/Plan`². The `./Plan` element has the required attributes `./@name` and `./majorFrame`; and contains a sequence of `./Slot` elements.

Each `./Slot` element has the following attributes:

`./Slot/@id` Slot Id's shall meet the id's rules defined in section 2.4. This value can be retrieved by the partition at run time, see section 5.7.1.

`./Slot/@duration` Time duration of the slot.

`./Slot/@partitionId` Id of the partition that will be executed during this slot.

`./Slot/@start` Offset with respect to the MAF start.

Slots intervals shall not overlap.

`./HwDescription/MemoryLayout` Defines the memory layout of the board. All the memory allocated to partitions, resident software and XtratuM itself shall be in the range of one of these areas.

`./HwDescription/Devices` The devices element contains the sequence the XtratuM devices. Currently XtratuM implements two types of devices: UART and memory blocks.

`./Uart` Has the required attributes `./Uart/@name`, `./Uart/@baudRate` and `./Uart/@id`. This element associates the hardware device `@id` with the `@name`, and programs the transmission speed.

`./MemoryBlockTable` This element contains a sequence of one or more `./Block` elements. A memory block device defines an area of RAM (ROM or FLASH) memory. This block of memory can then be used to store traces, health monitoring logs or the console output of a partition. Below is the list of attributes of the `./Block` element:

`./MemoryBlockTable/Block/@name` Required. Name which identifies the device. This name is only used to refer this device in the configuration file. Once compiled the configuration file this name is removed.

`./MemoryBlockTable/Block/@start` Required. Starting address of the memory block.

`./MemoryBlockTable/Block/@size` Required. Size of the memory block.

8.3.5 The `./SystemDescription/ResidentSw` element

The element `./PhysicalMemoryAreas` is used to declare the memory areas where the resident software will be located. This information is included in the configuration file for completeness (all the memory areas of the board shall be described in the configuration file) and used only to check memory overlaps errors.

Also the attribute `./@entryPoint` is used by XtratuM in the case of a cold system reset. In that case, XtratuM will give back the control of the system to the resident software by jumping to this address.

8.3.6 The `./SystemDescription/PartitionTable/Partition` element

Attribute description:

`./@id` Required. See the section 2.4 for a description on how to identify XtratuM objects.

`./@name` Optional.

²The large number of nested elements is for future compatibility with multiple plans and scheduling policies.

./@console Optional. The console device where the output of the hypercall `XM_write_console()` is copied to. 1635

./@flags Optional. List of features. Possible values are:

fp If set, the partition is allowed to use floating point operations. By default not set.

system If set, the partition has system privileges. By default not set.

Partition elements:

./PhysicalMemoryAreas Sequence of memory areas allocated to the partition. 1640

./HwResources Contains the list of interrupts and IO ports allocated to the partition.

./PortTable Contains the sequence of communication ports (queuing and sampling ports) of the partition.

./Trace Configuration of the trace facility of the partition. Same attributes than that of the `/SystemDescription/XMHypervisor/Trace` element. 1645

./TemporalRequirements An element which has two mandatory attributes: **./@period** and **./@duration**. **This data is not checked by XtratuM. Reserved for future use.**

Configuration of memory areas

The attributes are **@start**, **@size** and **@flags**. The **@flags** attribute is a list of the following values:

Value	Description
unmapped	Allocated to the partition, but not mapped by XtratuM in the page table.
shared	It is allowed to map this area in other partitions. 1650
read-only	The area is write-protected to the partition.
uncacheable	Memory cache is disabled.
rom	Not applicable in SPARC v8 boards. Only used in ia32 systems.

Configuration of I/O ports

There are two ways to allocate a port to a partition: using ranges of ports, and using the restricted port allocation. Both are declared by elements contained in the `./Partition/HwResources/IOPorts` element:

./Range A range of port addresses is allocated to the partition. The attributes of a range element are:

./Range/@base Required hexadecimal base address. 1655

./Range/@noPorts Required number of ports in this range. Each port is a word (4 bytes).

./Restricted An I/O port which is partially controlled by the partition. The attributes are:

./Restricted/@address Required hexadecimal address of the port.

./Restricted/@mask Optional (4 bytes hexadecimal). The bits set in this mask can be read and written by the partition. 1660

Those bits not allocated to this partition (i.e. the bit not set in the bitmask) can be allocated to other partitions.

Configuration of interrupts

The element `./Partition/HwResources/Interrupts` has the attribute `./@lines` which is a list of the interrupt number (in the range 0 to 16) allocated to the partition.

8.3.7 The `./SystemDescription/Channels` element

1665 This is an optional element with no attributes and which contains a list of channel elements. There are two types of channels:

`./SamplingChannel` Shall contain one `./Source` element and one or more `./Destination` elements. It has the following attributes:

1670 `./@maxLength` Required. The maximum message size that can be stored on this channel.

`./@refreshPeriod` Optional. The duration of validity of a written message. When a message is read after this period, the validity flag will be false.

`./QueuingChannel` Shall contain one `./Source` element and one `./Destination` element. It has the following attributes:

1675 `./@maxLength` Required. The maximum message size that can be stored on this channel.

`./@maxNoMessages` Required. The maximum number of messages that will be stored in the channel.

1680 **Note:** The `./QueuingChannel/@validPeriod` attribute has been removed with respect to XtratuM-2.2.x versions.

The arguments `maxNoMsgs` and `maxMsgSize` of the hypercalls `XM_create_queuing_port()` and `XM_create_sampling_port()` shall match the values of the attributes `./@maxNoMessages` and `./@maxNoMessages`.

The XML schema which defines the configuration file is in the appendix A.

Chapter 9

Tools

This section describes the tools to assist the integrator and the partition developers in the process of building the final system file. 1685

xmcparser: System XML configuration parser.

xmeformat: Converts ELF files into XEF ones.

xmpack: Creates the container file.

rswbuild: Creates a bootable file image.

9.1 XML configuration parser (**xmcparser**)

The utility `xmcparser` translates the XML configuration file containing the system description into binary form that can be directly used by XtratuM. 1690

In the first place, the configuration file is checked both, syntactically, and semantically (i.e. the data is correct). This tool uses the `libxml2` library to read, parse and validate the configuration file against the XML schema specification. Once validated by the library, the `xmcparser` performs a set of non-syntactical checks: 1695

- Memory area overlapping.
- Memory region overlapping.
- Memory area inside any region.
- Duplicated Partition's name and id.
- Allocated Cpus. 1700
- Replicated port's names and id.
- Cyclic scheduling plan.
- Cyclic scheduling plan slot partition ids.
- Hardware irqs allocated to partitions.
- Io port alignment. 1705
- Io ports allocated to partitions.
- Allowed health monitoring actions.

9.1.1 xmcparser

Compiles XtratuM XML configuration files

1710 SYNOPSIS

```
xmcparser [-c] [-s xsd_file] [-o output_file] XM.CF.xml
xmcparser -d
```

DESCRIPTION

1715 xmcparser reads an XtratuM XML configuration file and transforms it into a binary file which can be used directly by XtratuM at run time. xmcparser performs internally the following steps:

1. Parse the XML file.
2. Validate the XML data.
3. Generate a set of "C" data structures initialised with the XML data.
4. Compiles and links, using the target compiler, the "C" data structures. An ELF file is produced.
- 1720 5. The data section which contains the data in binary format is extracted and copied to the output file.

OPTIONS

-d

Prints the default XML schema used to validate the XML configuration file.

1725 -o file

Place output in file.

-s xsd_file

Use the XML schema xsd_file rather than the default XtratuM schema.

-c

1730 Stop after the stage of "C" generation; do not compile. The output is in the form of a "C" file.

9.2 ELF to XEF (xmeformat)

9.2.1 xmeformat

Creates and display information of XEF files

SYNOPSIS

```
xmeformat read [-h|-s|-m] file
1735 xmeformat build [-m] [-o outfile] [-c] [-p payload_file] file
```

DESCRIPTION

xmeformat converts an ELF, or a binary file, into an XEF format (XtratuM Executable Format). An XEF file contains one or more segments. A segment is a block of data that shall be copied in a contiguous area of memory (when loaded in main memory). The content of the XEF can optionally be compressed.

An XEF file has a header and a set of segments. The segments corresponds to the allocatable sections of the source ELF file. In the header, there is a reserved area (16 bytes) to store user defined information. This information is called user payload.

build

A new XEF file is created, using file as input.

-m

The source file is not an ELF file but a user defined customisation. In this case, no consistency checks are performed.

Customisation files are used to attach data to the partitions (See the xmpack command). This data will be accessible to the partition at boot time. It is commonly used as partition defined run-time configuration parameters.

-o file

Places output in file file.

-c

The XEF segments are compressed using the LSZZ algorithm.

-p file

The first 16 bytes of the file are copied into the payload area of the XEF header. The size of the file shall be at least 16 bytes, otherwise an error is returned.

The MD5 sum value is printed if no errors.

read

Shows the contents of the XEF file.

-h

Print the content of the header.

-s

Lists the segments and its attributes.

-m

Lists the table of custom files. This options only works for partition and hypervisor XEF files.

USAGE EXAMPLES

Create a customisation file:

```
$ xmeformat build -m -o custom_file.xef data.in
b07715208bbfe72897a259619e7d7a6d custom_file.xef
```

List the header of the XEF custom file:

```

$ xmeformat read -h custom_file.xef
XEF header:
  signature: 0x24584546
1775  version: 1.0.0
  flags: XEF_DIGEST XEF_CONTENT_CUSTOMFILE
  digest: b07715208bbfe72897a259619e7d7a6d
  payload: 00 00 00 00 00 00 00 00
          00 00 00 00 00 00 00 00
1780  file size: 232
  segment table offset: 80
  no. segments: 1
  customFile table offset: 104
  no. customFiles: 0
1785  image offset: 104
  image length: 127
  XM image's header: 0x0

```

Build the hypervisor XEF file:

```
$ xmeformat build -o xm_core.xef -c core/xm_core
```

```

1790  List the segments and headers of the XtratuM XEF file: $ xmeformat read -s xm_core.xef Segment
table: 1 segments segment 0 physical address: 0x40000000 virtual address: 0x40000000 file size:
68520 compressed file size: 32923 (48.05%)

```

```

$ xmeformat read -h xm_core.xef
XEF header:
1795  signature: 0x24584546
  version: 1.0.0
  flags: XEF_DIGEST XEF_COMPRESSED XEF_CONTENT_HYPERVISOR
  digest: 6698cfcf9311325e46e79ed50dfc9683
  payload: 00 00 00 00 00 00 00 00
1800          00 00 00 00 00 00 00 00
  file size: 33040
  segment table offset: 80
  no. segments: 1
  customFile table offset: 104
1805  no. customFiles: 1
  image offset: 112
  image length: 68520
  XM image's header: 0x40010b78
  compressed image length: 32928 (48.06%)

```

9.3 Container builder (xmpack)

1810 9.3.1 xmpack

Create an XtratuM system image container

SYNOPSIS

xmpack build -h xm_file[@offset]:conf_file[@offset] [-p id:part_file[@offset][:custom_file[@offset]]*] + container

xmpack list -c container

1815

DESCRIPTION

xmpack manipulates the XtratuM system container. The container is a simple filesystem designed to contain the XtratuM hypervisor core and zero or more XEF files. The container is an envelope to deploy all the system (hypervisor and partitions) from the host to the target. At boot time, the resident software is in charge of reading the contents of the container and coping the components to the RAM areas where the hypervisor and the partitions will be executed. Note that XtratuM has no knowledge about the container structure.

1820

The container is organised as a list of *components*. Each component is a list of XEF files. A component is used to store an executable unit, which can be: the XtratuM hypervisor or a partition. Each component is a list of one or more files. The first file shall be a valid XtratuM image (see the XtratuM binary file header) with the configuration file (once parsed and compiled into XEF format). The rest of the components are optional.

1825

xmpack is a helper utility that can be used to deploy an XtratuM system. It is not mandatory to use this tool to deploy the application (hypervisor and the partitions) in the target machine.

The following checks are done:

1830

- The binary image of the partitions fits into the allocated memory (as defined in the XM.CF).
- The size of the customisation files fits into the area reserved by each partition.
- The memory allocated to XtratuM is big enough to hold the XtratuM image plus the configuration file.

build

1835

A new *container* is created. Two kind of components can be defined:

-h to create an [H]ypervisor component:

The hypervisor entry is composed of the name of the XtratuM xef file and the binary configuration file (the result of processing the XM.CF file).

-p to create a [P]artition. The partition entries are composed of:

1840

The *id* of the partition, as specified in the XM.CF file. Note that this is the mechanism to bind the configuration description with the actual image of the partition. The part_file which shall contains the executable image. And zero or more custom_files. There shall be the same number of customisation files than that specified in the field `noCustomFiles` of the `xmImageHdr` structure.

1845

The elements that are part of each component are separated by ":".

By default, xmpack stores the files sequentially in the container. If the *offset* parameter is specified, then the file is placed at the given offset. The offset is defined with respect to the start of the container. The specified offset shall not overlap with existing data. The remaining files of the container will be placed after the end of this file.

1850

list

Shows the contents (components and the files of each component) of a container. If the option **-c** is given, the blocks allocated to each file are also shown.

USAGE EXAMPLES

1855 A new container with one hypervisor and one booting partition. The hypervisor container has two files: the hypervisor binary and the configuration table:

```
$ xmpack build build -h ../core/xm_core.bin:xm_ct.bin -p partition1.bin -o container
```

The same example but the second container has now two files: the partition image and a customisation file:

```
1860 $ xmpack/xmpack build -h ../core/xm_core.bin:xm_cf.bin \
    -p partition1.bin:p1.cfg \
    -p partition2.bin:p2.cfg container.bin
```

List the contents of the container:

```
$ xmpack list container.bin
1865 <Package file="container.bin" version="1.0.0">
    <XMHypervisor file="../core/xm_core.bin" fileSize="97188" offset="0x0" size="97192" >
        <Module file="xm_cf.bin" size="8976" />
    </XMHypervisor>
    <Partition file="partition1.bin" fileSize="29996" offset="0x19eb8" size="30000" >
1870     <Module file="p1.cfg" size="16" />
    </Partition>
    <Partition file="partition2.bin" fileSize="30292" offset="0x213f8" size="30296" >
        <Module file="p2.cfg" size="16" />
    </Partition>
1875 </Package>
```

9.4 Bootable image creator (rswbuild)

9.4.1 rswbuild

Create a bootable image

SYNOPSIS

```
rswbuild container bootable
```

1880 DESCRIPTION

rswbuild is a shell script that creates a bootable file by combining the resident software code with the container file. The container shall be a valid file created with the xmpack tool.

The resident software object file is read from the distribution directory pointer by the \$XTRATUM_PATH variable.

1885 USAGE EXAMPLES

```
rswbuild container resident_sw
```

Chapter 10

Security issues

This chapter introduces several security issues related with XtratuM which should be taken into account by partition developers.

10.1 Invoking a hypercall from libXM

Invoking a hypercall requires a non-standard protocol which must be directly implemented in assembly code. 1890

LibXM is a partition-level “C” library deployed jointly with XtratuM aiming to hide this complexity and ease the development of “C” partitions.

From the security point of view, XtratuM implements two stacks for each partition: one managed by the partition (user context) and another, internal, managed directly by XtratuM (supervisor context). 1895
The partition stack is used by the libXM to prepare the call to XtratuM (pretty much like the gLibc does). Once the hypercall service is invoked, XtratuM changes the stack to its own stack. This second stack may contain sensitive information, but it is located inside the memory space of XtratuM (not exposed). It is normal to observe that the partition stack is modified when a hypercall is called, however this behaviour is far from being considered an actual security issue. 1900

10.2 Preventing covert/side channels due to scheduling slot overrun

This version of XtratuM is non-preemptible: once the kernel starts an activity (e.g. a service), it cannot be interrupted until its completion. This behaviour includes any hypercall invocation: if a partition calls an hypercall just before a partition context switch must be performed, XtratuM will not carry out the action until the hypercall is finished. This overrun can be exploited to gain information. The information is obtained by measuring the temporal cost of the last hypercall. There are two types of information that can be retrieved: 1905

1. Whether the target partition was executing an hypercall at the end of the slot or not. If the spy partition start at the nominal slot start time or not. 1910
2. In the case of being executing a hypercall, how much time XtratuM needed to attend it: the cost of the last hypercall.

In the case of a covert channel, the maximum bandwidth is determined by the duration of the longest hypercall divided by the clock resolution. A rough estimation (supposing that the maximum message length is 4096 Bytes) is 4 bits at each partition context switch.

In the case of a side channel, the bandwidth is drastically reduced due to the uncertainty/randomness introduced by the execution of the target partition.

There are several strategies to address this issue:

1. At integrator level: design a scheduling plan that lest some idle time before the end of a slot and the beginning of the next one. The Xoncrete scheduling tool is able to implement that solution automatically. Figure 10.1 shows two scenarios: scenario 1 where the integrator has not left spare time between one partition slot and the next one, enabling the partition in light grey overrunning the start of the dark gray one. Scenario 2 sketches the same case but leaving spare time between one slot and the next one. So, in this case, execution overruns can not occur.
2. At partition level: stop invoking hypercalls some time before the end of the slot. This way, there will be no hypercalls being executed when the slot end occurs, so the next partition will start always with no delay.
3. At hypervisor level:
 - (a) Change the design of XtratuM to convert it preemptable. Hypercalls would be interrupted when the end of the slot is reached, and later resumed when the partition is active again.
 - (b) Implement the partial preemptability in XtratuM (voluntary preemption). XtratuM is by default atomic (non-preemptable) but at some designated safe places in the code, the preemption is allowed.

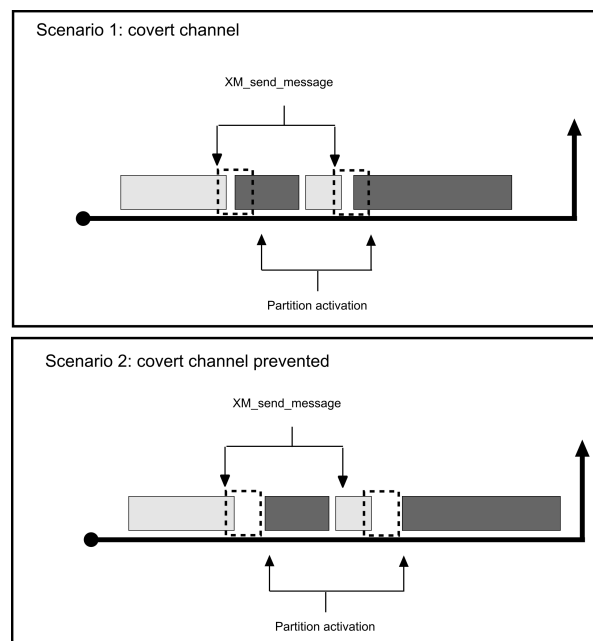


Figure 10.1: Covert channel caused by an incorrect scheduling plan and a solution.

Appendix A

XML Schema Definition

A.1 XML Schema file

basicstyle

```
1 <?xml version="1.0"?>
2 <xs:schema targetNamespace="http://www.xtratum.org/xm-3.x"
3           xmlns:xs="http://www.w3.org/2001/XMLSchema"
4           xmlns="http://www.xtratum.org/xm-3.x"
5           elementFormDefault="qualified"
6           attributeFormDefault="unqualified">
7
8   <!-- Basic types definition -->
9   <xs:simpleType name="id_t">
10     <xs:restriction base="xs:integer">
11       <xs:minInclusive value="0"/>
12     </xs:restriction>
13   </xs:simpleType>
14
15   <xs:simpleType name="idString_t">
16     <xs:restriction base="xs:string">
17       <xs:minLength value="1"/>
18     </xs:restriction>
19   </xs:simpleType>
20
21   <xs:simpleType name="hwIrqId_t">
22     <xs:restriction base="xs:integer">
23       <xs:minInclusive value="0"/>
24       <xs:maxExclusive value="16"/>
25     </xs:restriction>
26   </xs:simpleType>
27
28   <xs:simpleType name="hwIrqIdList_t">
29     <xs:list itemType="hwIrqId_t"/>
30   </xs:simpleType>
31
32   <xs:simpleType name="hex_t">
33     <xs:restriction base="xs:string">
34       <xs:pattern value="0x[0-9a-fA-F]+"/>

```

```

35     </xs:restriction>
36 </xs:simpleType>
37 <xs:simpleType name="version_t">
38     <xs:restriction base="xs:string">
39         <xs:pattern value="[0-9]+.[0-9]+.[0-9]+" />
40     </xs:restriction>
41 </xs:simpleType>
42
43 <xs:simpleType name="freqUnit_t">
44     <xs:restriction base="xs:string">
45         <xs:pattern value="[0-9]+(.[0-9]+)?([MK][Hh]z)" />
46     </xs:restriction>
47 </xs:simpleType>
48
49 <xs:simpleType name="processorFeatures_t">
50     <xs:restriction base="xs:string">
51         <xs:enumeration value="XM_CPU_LEON2_WA1" />
52     </xs:restriction>
53 </xs:simpleType>
54
55 <xs:simpleType name="processorFeaturesList_t">
56     <xs:list itemType="processorFeatures_t" />
57 </xs:simpleType>
58
59 <xs:simpleType name="partitionFlags_t">
60     <xs:restriction base="xs:string">
61         <xs:enumeration value="system" />
62         <xs:enumeration value="fp" />
63     </xs:restriction>
64 </xs:simpleType>
65
66 <xs:simpleType name="partitionFlagsList_t">
67     <xs:list itemType="partitionFlags_t" />
68 </xs:simpleType>
69
70 <xs:simpleType name="sizeUnit_t">
71     <xs:restriction base="xs:string">
72         <xs:pattern value="[0-9]+(.[0-9]+)?([MK]?B)" />
73     </xs:restriction>
74 </xs:simpleType>
75
76 <xs:simpleType name="timeUnit_t">
77     <xs:restriction base="xs:string">
78         <xs:pattern value="[0-9]+(.[0-9]+)?([\mu]?[sS])" />
79     </xs:restriction>
80 </xs:simpleType>
81
82 <xs:simpleType name="hmString_t">
83     <xs:restriction base="xs:string">
84         <xs:enumeration value="XM_HM_EV_INTERNAL_ERROR" />
85         <xs:enumeration value="XM_HM_EV_UNEXPECTED_TRAP" />
86         <xs:enumeration value="XM_HM_EV_PARTITION_ERROR" />
87         <xs:enumeration value="XM_HM_EV_PARTITION_INTEGRITY" />
88         <xs:enumeration value="XM_HM_EV_MEM_PROTECTION" />

```

```
89     <xs:enumeration value="XM_HM_EV_OVERRUN"/>
90     <xs:enumeration value="XM_HM_EV_SCHED_ERROR"/>
91     <xs:enumeration value="XM_HM_EV_WATCHDOG_TIMER"/>
92     <xs:enumeration value="XM_HM_EV_INCOMPATIBLE_INTERFACE"/>
93     <xs:enumeration value="XM_HM_EV_WRITE_ERROR"/>
94     <xs:enumeration value="XM_HM_EV_INSTR_ACCESS_MMU_MISS"/>
95     <xs:enumeration value="XM_HM_EV_INSTR_ACCESS_ERROR"/>
96     <xs:enumeration value="XM_HM_EV_UNIMPLEMENTED_FLUSH"/>
97     <xs:enumeration value="XM_HM_EV_WATCHPOINT_DETECTED"/>
98     <xs:enumeration value="XM_HM_EV_DATA_ACCESS_ERROR"/>
99     <xs:enumeration value="XM_HM_EV_DATA_ACCESS_MMU_MISS"/>
100    <xs:enumeration value="XM_HM_EV_INSTR_ACCESS_EXCEPTION"/>
101    <xs:enumeration value="XM_HM_EV_ILLEGAL_INSTR"/>
102    <xs:enumeration value="XM_HM_EV_PRIVILEGED_INSTR"/>
103    <xs:enumeration value="XM_HM_EV_FP_DISABLED"/>
104    <xs:enumeration value="XM_HM_EV_CP_DISABLED"/>
105    <xs:enumeration value="XM_HM_EV_REGISTER_HARDWARE_ERROR"/>
106    <xs:enumeration value="XM_HM_EV_MEM_ADDR_NOT_ALIGNED"/>
107    <xs:enumeration value="XM_HM_EV_FP_EXCEPTION"/>
108    <xs:enumeration value="XM_HM_EV_CP_EXCEPTION"/>
109    <xs:enumeration value="XM_HM_EV_DATA_ACCESS_EXCEPTION"/>
110    <xs:enumeration value="XM_HM_EV_TAG_OVERFLOW"/>
111    <xs:enumeration value="XM_HM_EV_DIVIDE_EXCEPTION"/>
112  </xs:restriction>
113 </xs:simpleType>
114
115
116 <xs:simpleType name="hmAction_t">
117   <xs:restriction base="xs:string">
118     <xs:enumeration value="XM_HM_AC_IGNORE"/>
119     <xs:enumeration value="XM_HM_AC_SHUTDOWN"/>
120     <xs:enumeration value="XM_HM_AC_COLD_RESET"/>
121     <xs:enumeration value="XM_HM_AC_WARM_RESET"/>
122     <xs:enumeration value="XM_HM_AC_SUSPEND"/>
123     <xs:enumeration value="XM_HM_AC_HALT"/>
124     <xs:enumeration value="XM_HM_AC_PROPAGATE"/>
125     <xs:enumeration value="XM_HM_AC_SWITCH_TO_MAINTENANCE" />
126   </xs:restriction>
127 </xs:simpleType>
128
129 <xs:simpleType name="memAreaFlags_t">
130   <xs:restriction base="xs:string">
131     <xs:enumeration value="unmapped"/>
132     <xs:enumeration value="shared"/>
133     <xs:enumeration value="read-only"/>
134     <xs:enumeration value="uncacheable"/>
135     <xs:enumeration value="rom"/>
136     <xs:enumeration value="shlib"/>
137     <xs:enumeration value="flag0"/>
138     <xs:enumeration value="flag1"/>
139     <xs:enumeration value="flag2"/>
140     <xs:enumeration value="flag3"/>
141   </xs:restriction>
142 </xs:simpleType>
```

```

143
144 <xs:simpleType name="memAreaFlagsList_t">
145   <xs:list itemType="memAreaFlags_t"/>
146 </xs:simpleType>
147
148 <xs:simpleType name="memRegion_t">
149   <xs:restriction base="xs:string">
150     <xs:enumeration value="sdram"/>
151     <xs:enumeration value="stram"/>
152   </xs:restriction>
153 </xs:simpleType>
154
155 <xs:simpleType name="portType_t">
156   <xs:restriction base="xs:string">
157     <xs:enumeration value="queuing"/>
158     <xs:enumeration value="sampling"/>
159   </xs:restriction>
160 </xs:simpleType>
161
162 <xs:simpleType name="direction_t">
163   <xs:restriction base="xs:string">
164     <xs:enumeration value="source"/>
165     <xs:enumeration value="destination"/>
166   </xs:restriction>
167 </xs:simpleType>
168
169 <xs:simpleType name="yntf_t">
170   <xs:restriction base="xs:string">
171     <xs:enumeration value="yes"/>
172     <xs:enumeration value="no"/>
173     <xs:enumeration value="true"/>
174     <xs:enumeration value="false"/>
175   </xs:restriction>
176 </xs:simpleType>
177 <!-- End Types -->
178
179 <!-- Elements -->
180 <!-- Hypervisor -->
181 <xs:complexType name="hypervisor_e">
182   <xs:all>
183     <xs:element name="PhysicalMemoryAreas" type="memoryArea_e"/>
184     <xs:element name="HealthMonitor" type="healthMonitor_e" minOccurs="0" />
185     <xs:element name="Trace" type="trace_e" minOccurs="0" />
186   </xs:all>
187   <xs:attribute name="console" type="idString_t" use="optional"/>
188   <xs:attribute name="healthMonitorDevice" type="idString_t" use="optional"/>
189 </xs:complexType>
190
191 <!-- Rsw -->
192 <xs:complexType name="rsw_e">
193   <xs:all>
194     <xs:element name="PhysicalMemoryAreas" type="memoryArea_e"/>
195   </xs:all>
196 </xs:complexType>

```



```
197
198 <!-- Partition -->
199 <xs:complexType name="partition_e">
200   <xs:all>
201     <xs:element name="PhysicalMemoryAreas" type="memoryArea_e"/>
202     <xs:element name="TemporalRequirements" minOccurs="0">
203       <xs:complexType>
204         <xs:attribute name="period" type="timeUnit_t" use="required"/>
205         <xs:attribute name="duration" type="timeUnit_t" use="required"/>
206       </xs:complexType>
207     </xs:element>
208     <xs:element name="HealthMonitor" type="healthMonitor_e" minOccurs="0" />
209     <xs:element name="HwResources" type="hwResources_e" minOccurs="0" />
210     <xs:element name="PortTable" type="partitionPorts_e" minOccurs="0" />
211     <xs:element name="Trace" type="trace_e" minOccurs="0" />
212   </xs:all>
213   <xs:attribute name="id" type="id_t" use="required"/>
214   <xs:attribute name="name" type="idString_t" use="optional"/>
215   <xs:attribute name="console" type="idString_t" use="optional"/>
216   <xs:attribute name="flags" type="partitionFlagsList_t" use="optional" />
217 </xs:complexType>
218
219 <!-- Trace -->
220 <xs:complexType name="trace_e">
221   <xs:attribute name="device" type="idString_t" use="required"/>
222   <xs:attribute name="bitmask" type="hex_t" use="required"/>
223 </xs:complexType>
224
225 <!-- Communication Ports -->
226 <xs:complexType name="partitionPorts_e">
227   <xs:sequence minOccurs="0" maxOccurs="unbounded">
228     <xs:element name="Port">
229       <xs:complexType>
230         <xs:attribute name="name" type="idString_t" use="required"/>
231         <xs:attribute name="direction" type="direction_t" use="required"/>
232         <xs:attribute name="type" type="portType_t" use="required"/>
233       </xs:complexType>
234     </xs:element>
235   </xs:sequence>
236 </xs:complexType>
237
238 <!-- Channels -->
239 <xs:complexType name="channels_e">
240   <xs:sequence minOccurs="0" maxOccurs="unbounded">
241     <xs:choice>
242       <xs:element name="SamplingChannel">
243         <xs:complexType>
244           <xs:sequence minOccurs="1">
245             <xs:element name="Source" type="ipcPort_e" />
246             <xs:sequence minOccurs="1" maxOccurs="unbounded">
247               <xs:element name="Destination" type="ipcPort_e"/>
248             </xs:sequence>
249           </xs:sequence>

```

```

250         <xs:attribute name="maxLength" type="sizeUnit_t" use="
                required"/>
251         <xs:attribute name="refreshPeriod" type="timeUnit_t" use="optional"/
                >
252     </xs:complexType>
253 </xs:element>
254 <xs:element name="QueuingChannel">
255     <xs:complexType>
256         <xs:all minOccurs="1">
257             <xs:element name="Source" type="ipcPort_e" />
258             <xs:element name="Destination" type="ipcPort_e"/>
259         </xs:all>
260         <xs:attribute name="maxLength" type="sizeUnit_t" use="
                required"/>
261         <xs:attribute name="maxNoMessages" type="xs:positiveInteger" use="
                required"/>
262     </xs:complexType>
263 </xs:element>
264 </xs:choice>
265 </xs:sequence>
266 </xs:complexType>
267
268 <!-- Devices -->
269 <xs:complexType name="devices_e">
270     <xs:sequence minOccurs="0" maxOccurs="unbounded">
271         <xs:choice>
272
273             <xs:element name="MemoryBlock" minOccurs="0">
274                 <xs:complexType>
275                     <xs:attribute name="name" type="idString_t" use="required"/>
276                     <xs:attribute name="start" type="hex_t" use="required"/>
277                     <xs:attribute name="size" type="sizeUnit_t" use="required"/>
278                 </xs:complexType>
279             </xs:element>
280
281
282             <xs:element name="Uart" minOccurs="0">
283                 <xs:complexType>
284                     <xs:attribute name="name" type="idString_t" use="required"/>
285                     <xs:attribute name="id" type="idString_t" use="required"/>
286                     <xs:attribute name="baudRate" type="xs:positiveInteger" use="
                            required"/>
287                     <xs:attribute name="rxBufferLen" type="xs:positiveInteger" use="
                            optional" />
288                 </xs:complexType>
289             </xs:element>
290
291         </xs:choice>
292     </xs:sequence>
293 </xs:complexType>
294
295 <!-- IPC Port -->
296 <xs:complexType name="ipcPort_e">
297     <xs:attribute name="partitionId" type="id_t" use="required"/>

```

```
298     <xs:attribute name="partitionName" type="idString_t" use="optional"/>
299     <xs:attribute name="portName" type="idString_t" use="required"/>
300 </xs:complexType>
301
302 <!-- Hw Description -->
303 <xs:complexType name="hwDescription_e">
304     <xs:all>
305         <xs:element name="ProcessorTable">
306             <xs:complexType>
307                 <xs:sequence minOccurs="1" maxOccurs="256">
308                     <xs:element name="Processor" type="processor_e" />
309                 </xs:sequence>
310             </xs:complexType>
311         </xs:element>
312         <xs:element name="MemoryLayout" type="memoryLayout_e"/>
313         <xs:element name="Devices" type="devices_e"/>
314     </xs:all>
315 </xs:complexType>
316
317 <!-- Processor -->
318 <xs:complexType name="processor_e">
319     <xs:all>
320         <xs:element name="CyclicPlanTable" type="cyclicPlan_e"/>
321     </xs:all>
322     <xs:attribute name="id" type="id_t" use="required"/>
323     <xs:attribute name="frequency" type="freqUnit_t" use="optional"/>
324     <xs:attribute name="features" type="processorFeaturesList_t" use="optional"
325     />
326 </xs:complexType>
327
328 <!-- HwResource -->
329 <xs:complexType name="hwResources_e">
330     <xs:all>
331         <xs:element name="IoPorts" type="ioPorts_e" minOccurs="0" />
332         <xs:element name="Interrupts" minOccurs="0">
333             <xs:complexType>
334                 <xs:attribute name="lines" type="hwIrqIdList_t" use="required"/>
335             </xs:complexType>
336         </xs:element>
337     </xs:all>
338 </xs:complexType>
339
340 <!-- Io Ports -->
341 <xs:complexType name="ioPorts_e">
342     <xs:sequence minOccurs="0" maxOccurs="unbounded">
343         <xs:choice>
344             <xs:element name="Range">
345                 <xs:complexType>
346                     <xs:attribute name="base" type="hex_t" use="required"/>
347                     <xs:attribute name="noPorts" type="xs:positiveInteger" use="required"
348                     "/>
349                 </xs:complexType>
350             </xs:element>
351             <xs:element name="Restricted">
```

```

350         <xs:complexType>
351             <xs:attribute name="address" type="hex_t" use="required"/>
352             <xs:attribute name="mask" type="hex_t" use="optional"/>
353         </xs:complexType>
354     </xs:element>
355 </xs:choice>
356 </xs:sequence>
357 </xs:complexType>
358
359 <!-- CyclicPlan -->
360 <xs:complexType name="cyclicPlan_e">
361     <xs:sequence minOccurs="1" maxOccurs="unbounded">
362         <xs:element name="Plan" type="plan_e" />
363     </xs:sequence>
364 </xs:complexType>
365
366 <!-- Plan -->
367 <xs:complexType name="plan_e">
368     <xs:sequence minOccurs="1" maxOccurs="unbounded">
369         <xs:element name="Slot">
370             <xs:complexType>
371                 <xs:attribute name="id" type="id_t" use="required"/>
372                 <xs:attribute name="start" type="timeUnit_t" use="required"/>
373                 <xs:attribute name="duration" type="timeUnit_t" use="required"/>
374                 <xs:attribute name="partitionId" type="id_t" use="required"/>
375             </xs:complexType>
376         </xs:element>
377     </xs:sequence>
378     <xs:attribute name="id" type="id_t" use="required"/>
379     <xs:attribute name="majorFrame" type="timeUnit_t" use="required"/>
380 </xs:complexType>
381
382 <!-- Health Monitor -->
383 <xs:complexType name="healthMonitor_e">
384     <xs:sequence minOccurs="1" maxOccurs="unbounded">
385         <xs:element name="Event">
386             <xs:complexType>
387                 <xs:attribute name="name" type="hmString_t" use="required"/>
388                 <xs:attribute name="action" type="hmAction_t" use="required"/>
389                 <xs:attribute name="log" type="yntf_t" use="required"/>
390             </xs:complexType>
391         </xs:element>
392     </xs:sequence>
393 </xs:complexType>
394
395 <!-- Memory Layout -->
396 <xs:complexType name="memoryLayout_e">
397     <xs:sequence minOccurs="1" maxOccurs="unbounded">
398         <xs:element name="Region">
399             <xs:complexType>
400                 <xs:attribute name="type" type="memRegion_t" use="required"/>
401                 <xs:attribute name="start" type="hex_t" use="required"/>
402                 <xs:attribute name="size" type="sizeUnit_t" use="required"/>
403             </xs:complexType>

```

```

404     </xs:element>
405     </xs:sequence>
406 </xs:complexType>
407
408 <!-- Memory Area -->
409 <xs:complexType name="memoryArea_e">
410   <xs:sequence minOccurs="1" maxOccurs="unbounded">
411     <xs:element name="Area">
412       <xs:complexType>
413         <xs:attribute name="start" type="hex_t" use="required"/>
414         <xs:attribute name="size" type="sizeUnit_t" use="required"/>
415         <xs:attribute name="flags" type="memAreaFlagsList_t" use="optional"/>
416         <xs:attribute name="mappedAt" type="hex_t" use="optional" />
417       </xs:complexType>
418     </xs:element>
419   </xs:sequence>
420 </xs:complexType>
421
422 <!-- Root Element -->
423 <xs:element name="SystemDescription">
424   <xs:complexType>
425     <xs:all>
426       <xs:element name="HwDescription" type="hwDescription_e" />
427       <xs:element name="XMHypervisor" type="hypervisor_e"/>
428       <xs:element name="ResidentSw" type="rsw_e" minOccurs="0"/>
429       <xs:element name="PartitionTable">
430         <xs:complexType>
431           <xs:sequence maxOccurs="unbounded">
432             <xs:element name="Partition" type="partition_e" />
433           </xs:sequence>
434         </xs:complexType>
435       </xs:element>
436       <xs:element name="Channels" type="channels_e" minOccurs="0" />
437     </xs:all>
438     <xs:attribute name="version" type="version_t" use="required"/>
439     <xs:attribute name="name" type="idString_t" use="required"/>
440   </xs:complexType>
441 </xs:element>
442 <!-- End Root Element -->
443 <!-- Elements -->
444 </xs:schema>

```

Listing A.1: xmc.xsd

A.2 Configuration file example

```

<SystemDescription xmlns="http://www.xtratum.org/xm-3.x"
  version="1.0.0" name="hello_world">
  <HwDescription>
    <ProcessorTable>
      <Processor id="0" frequency="50Mhz">
        <Sched>
          <CyclicPlan>

```

```

        <Plan majorFrame="2ms">
            <Slot id="0" start="0ms" duration="1ms" partitionId="0"/>
            <Slot id="1" start="1ms" duration="1ms" partitionId="1"/>
        </Plan>
    </CyclicPlan>
</Sched>
</Processor>
</ProcessorTable>
<Devices>
    <Uart id="0" baudRate="115200" name="Uart" />
    <MemoryBlock name="MemDisk0" start="0x40100000" size="256KB" />
    <MemoryBlock name="MemDisk1" start="0x40150000" size="256KB" />
    <MemoryBlock name="MemDisk2" start="0x40200000" size="256KB" />
</Devices>
<MemoryLayout>
    <Region type="stam" start="0x40000000" size="4MB"/>
</MemoryLayout>
</HwDescription>

<XMHypervisor console="Uart">
    <PhysicalMemoryAreas>
        <Area start="0x40000000" size="512KB" flags="uncacheable"/>
    </PhysicalMemoryAreas>
    <HealthMonitor>
        <Event name="XM_HM_EV_INTERNAL_ERROR" action="XM_HM_AC_IGNORE" log="yes"/>
    </HealthMonitor>
    <Trace device="MemDisk0" bitmask="0xabcd" />
</XMHypervisor>

<ResidentSw>
    <PhysicalMemoryAreas>
        <Area start="0x40200000" size="1MB" flags="shared"/>
    </PhysicalMemoryAreas>
</ResidentSw>

<PartitionTable>
    <Partition id="0" name="Partition1" flags="system" console="Uart">
        <PhysicalMemoryAreas>
            <Area start="0x40080000" size="512KB" />
            <Area start="0x40200000" size="1MB" flags="shared"/>
        </PhysicalMemoryAreas>
        <TemporalRequirements duration="500ms" period="500ms"/>
        <HwResources>
            <IoPorts>
                <Restricted address="0xfc" mask="0xff"/>
                <Range base="0x80" noPorts="10"/>
            </IoPorts>
        </HwResources>
        <PortTable>
            <Port name="writerQ" type="queuing" direction="source" />
            <Port name="writerS" type="sampling" direction="source" />
        </PortTable>
    </Partition>
    <Partition id="1" name="Partition2" flags="system" console="Uart">
        <PhysicalMemoryAreas>
            <Area start="0x40100000" size="512KB" flags="uncacheable" />

```

```
</PhysicalMemoryAreas>
<TemporalRequirements duration="500ms" period="500ms"/>
<PortTable>
  <Port name="readerQ" type="queuing" direction="destination" />
  <Port name="readerS" type="sampling" direction="destination" />
</PortTable>
<HwResources>
  <Interrupts lines="4 5" />
  <IoPorts>
    <Restricted address="0x80000240" mask="0xff"/>
    <Range base="0x380" noPorts="10"/>
  </IoPorts>
</HwResources>
</Partition>
</PartitionTable>

<Channels>
  <QueuingChannel maxMessageLength="512B" maxNoMessages="10">
    <Source partitionId="0" portName="writerQ" />
    <Destination partitionId="1" portName="readerQ" />
  </QueuingChannel>
  <SamplingChannel maxMessageLength="512B">
    <Source partitionId="0" portName="writerS" />
    <Destination partitionId="1" portName="readerS" />
  </SamplingChannel>
</Channels>

</SystemDescription>
```

Listing A.2: /user/tools/xmcpaser/xm_cf.sparcv8.xml

This page is intentionally left blank.

GNU Free Documentation License

1935

Version 1.2, November 2002

Copyright © 2000,2001,2002 Free Software Foundation, Inc.

51 Franklin St, Fifth Floor, Boston, MA 02110-1301 USA

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

1940

Preamble

The purpose of this License is to make a manual, textbook, or other functional and useful document “free” in the sense of freedom: to assure everyone the effective freedom to copy and redistribute it, with or without modifying it, either commercially or noncommercially. Secondly, this License preserves for the author and publisher a way to get credit for their work, while not being considered responsible for modifications made by others.

1945

This License is a kind of “copyleft”, which means that derivative works of the document must themselves be free in the same sense. It complements the GNU General Public License, which is a copyleft license designed for free software.

We have designed this License in order to use it for manuals for free software, because free software needs free documentation: a free program should come with manuals providing the same freedoms that the software does. But this License is not limited to software manuals; it can be used for any textual work, regardless of subject matter or whether it is published as a printed book. We recommend this License principally for works whose purpose is instruction or reference.

1950

1. APPLICABILITY AND DEFINITIONS

1955

This License applies to any manual or other work, in any medium, that contains a notice placed by the copyright holder saying it can be distributed under the terms of this License. Such a notice grants a world-wide, royalty-free license, unlimited in duration, to use that work under the conditions stated herein. The “**Document**”, below, refers to any such manual or work. Any member of the public is a licensee, and is addressed as “**you**”. You accept the license if you copy, modify or distribute the work in a way requiring permission under copyright law.

1960

A “**Modified Version**” of the Document means any work containing the Document or a portion of it, either copied verbatim, or with modifications and/or translated into another language.

A “**Secondary Section**” is a named appendix or a front-matter section of the Document that deals exclusively with the relationship of the publishers or authors of the Document to the Document’s overall subject (or to related matters) and contains nothing that could fall directly within that overall subject. (Thus, if the Document is in part a textbook of mathematics, a Secondary Section may not explain

1965

any mathematics.) The relationship could be a matter of historical connection with the subject or with related matters, or of legal, commercial, philosophical, ethical or political position regarding them.

1970 The “**Invariant Sections**” are certain Secondary Sections whose titles are designated, as being those of Invariant Sections, in the notice that says that the Document is released under this License. If a section does not fit the above definition of Secondary then it is not allowed to be designated as Invariant. The Document may contain zero Invariant Sections. If the Document does not identify any Invariant Sections then there are none.

1975 The “**Cover Texts**” are certain short passages of text that are listed, as Front-Cover Texts or Back-Cover Texts, in the notice that says that the Document is released under this License. A Front-Cover Text may be at most 5 words, and a Back-Cover Text may be at most 25 words.

1980 A “**Transparent**” copy of the Document means a machine-readable copy, represented in a format whose specification is available to the general public, that is suitable for revising the document straightforwardly with generic text editors or (for images composed of pixels) generic paint programs or (for drawings) some widely available drawing editor, and that is suitable for input to text formatters or for automatic translation to a variety of formats suitable for input to text formatters. A copy made in an otherwise Transparent file format whose markup, or absence of markup, has been arranged to thwart or discourage subsequent modification by readers is not Transparent. An image format is not Transparent if used for any substantial amount of text. A copy that is not “Transparent” is called “**Opaque**”.

1985 Examples of suitable formats for Transparent copies include plain ASCII without markup, Texinfo input format, LaTeX input format, SGML or XML using a publicly available DTD, and standard-conforming simple HTML, PostScript or PDF designed for human modification. Examples of transparent image formats include PNG, XCF and JPG. Opaque formats include proprietary formats that can be read and edited only by proprietary word processors, SGML or XML for which the DTD and/or processing tools are not generally available, and the machine-generated HTML, PostScript or PDF produced by some word processors for output purposes only.

1995 The “**Title Page**” means, for a printed book, the title page itself, plus such following pages as are needed to hold, legibly, the material this License requires to appear in the title page. For works in formats which do not have any title page as such, “Title Page” means the text near the most prominent appearance of the work’s title, preceding the beginning of the body of the text.

2000 A section “**Entitled XYZ**” means a named subunit of the Document whose title either is precisely XYZ or contains XYZ in parentheses following text that translates XYZ in another language. (Here XYZ stands for a specific section name mentioned below, such as “**Acknowledgements**”, “**Dedications**”, “**Endorsements**”, or “**History**”.) To “**Preserve the Title**” of such a section when you modify the Document means that it remains a section “Entitled XYZ” according to this definition.

2005 The Document may include Warranty Disclaimers next to the notice which states that this License applies to the Document. These Warranty Disclaimers are considered to be included by reference in this License, but only as regards disclaiming warranties: any other implication that these Warranty Disclaimers may have is void and has no effect on the meaning of this License.

2. VERBATIM COPYING

2010 You may copy and distribute the Document in any medium, either commercially or noncommercially, provided that this License, the copyright notices, and the license notice saying this License applies to the Document are reproduced in all copies, and that you add no other conditions whatsoever to those of this License. You may not use technical measures to obstruct or control the reading or further copying of the copies you make or distribute. However, you may accept compensation in exchange for copies. If you distribute a large enough number of copies you must also follow the conditions in section 3.

You may also lend copies, under the same conditions stated above, and you may publicly display copies.

3. COPYING IN QUANTITY

2015

If you publish printed copies (or copies in media that commonly have printed covers) of the Document, numbering more than 100, and the Document's license notice requires Cover Texts, you must enclose the copies in covers that carry, clearly and legibly, all these Cover Texts: Front-Cover Texts on the front cover, and Back-Cover Texts on the back cover. Both covers must also clearly and legibly identify you as the publisher of these copies. The front cover must present the full title with all words of the title equally prominent and visible. You may add other material on the covers in addition. Copying with changes limited to the covers, as long as they preserve the title of the Document and satisfy these conditions, can be treated as verbatim copying in other respects.

2020

If the required texts for either cover are too voluminous to fit legibly, you should put the first ones listed (as many as fit reasonably) on the actual cover, and continue the rest onto adjacent pages.

2025

If you publish or distribute Opaque copies of the Document numbering more than 100, you must either include a machine-readable Transparent copy along with each Opaque copy, or state in or with each Opaque copy a computer-network location from which the general network-using public has access to download using public-standard network protocols a complete Transparent copy of the Document, free of added material. If you use the latter option, you must take reasonably prudent steps, when you begin distribution of Opaque copies in quantity, to ensure that this Transparent copy will remain thus accessible at the stated location until at least one year after the last time you distribute an Opaque copy (directly or through your agents or retailers) of that edition to the public.

2030

It is requested, but not required, that you contact the authors of the Document well before redistributing any large number of copies, to give them a chance to provide you with an updated version of the Document.

2035

4. MODIFICATIONS

You may copy and distribute a Modified Version of the Document under the conditions of sections 2 and 3 above, provided that you release the Modified Version under precisely this License, with the Modified Version filling the role of the Document, thus licensing distribution and modification of the Modified Version to whoever possesses a copy of it. In addition, you must do these things in the Modified Version:

2040

- A. Use in the Title Page (and on the covers, if any) a title distinct from that of the Document, and from those of previous versions (which should, if there were any, be listed in the History section of the Document). You may use the same title as a previous version if the original publisher of that version gives permission.
- 2045
- B. List on the Title Page, as authors, one or more persons or entities responsible for authorship of the modifications in the Modified Version, together with at least five of the principal authors of the Document (all of its principal authors, if it has fewer than five), unless they release you from this requirement.
- 2050
- C. State on the Title page the name of the publisher of the Modified Version, as the publisher.
 - D. Preserve all the copyright notices of the Document.
 - E. Add an appropriate copyright notice for your modifications adjacent to the other copyright notices.
 - F. Include, immediately after the copyright notices, a license notice giving the public permission to use the Modified Version under the terms of this License, in the form shown in the Addendum below.
- 2055

G. Preserve in that license notice the full lists of Invariant Sections and required Cover Texts given in the Document's license notice.

2060 H. Include an unaltered copy of this License.

I. Preserve the section Entitled "History", Preserve its Title, and add to it an item stating at least the title, year, new authors, and publisher of the Modified Version as given on the Title Page. If there is no section Entitled "History" in the Document, create one stating the title, year, authors, and publisher of the Document as given on its Title Page, then add an item describing the Modified Version as stated in the previous sentence.

J. Preserve the network location, if any, given in the Document for public access to a Transparent copy of the Document, and likewise the network locations given in the Document for previous versions it was based on. These may be placed in the "History" section. You may omit a network location for a work that was published at least four years before the Document itself, or if the original publisher of the version it refers to gives permission.

K. For any section Entitled "Acknowledgements" or "Dedications", Preserve the Title of the section, and preserve in the section all the substance and tone of each of the contributor acknowledgements and/or dedications given therein.

L. Preserve all the Invariant Sections of the Document, unaltered in their text and in their titles. Section numbers or the equivalent are not considered part of the section titles.

M. Delete any section Entitled "Endorsements". Such a section may not be included in the Modified Version.

N. Do not retitle any existing section to be Entitled "Endorsements" or to conflict in title with any Invariant Section.

2080 O. Preserve any Warranty Disclaimers.

If the Modified Version includes new front-matter sections or appendices that qualify as Secondary Sections and contain no material copied from the Document, you may at your option designate some or all of these sections as invariant. To do this, add their titles to the list of Invariant Sections in the Modified Version's license notice. These titles must be distinct from any other section titles.

2085 You may add a section Entitled "Endorsements", provided it contains nothing but endorsements of your Modified Version by various parties—for example, statements of peer review or that the text has been approved by an organization as the authoritative definition of a standard.

You may add a passage of up to five words as a Front-Cover Text, and a passage of up to 25 words as a Back-Cover Text, to the end of the list of Cover Texts in the Modified Version. Only one passage of Front-Cover Text and one of Back-Cover Text may be added by (or through arrangements made by) any one entity. If the Document already includes a cover text for the same cover, previously added by you or by arrangement made by the same entity you are acting on behalf of, you may not add another; but you may replace the old one, on explicit permission from the previous publisher that added the old one.

2095 The author(s) and publisher(s) of the Document do not by this License give permission to use their names for publicity for or to assert or imply endorsement of any Modified Version.

5. COMBINING DOCUMENTS

You may combine the Document with other documents released under this License, under the terms defined in section 4 above for modified versions, provided that you include in the combination all of the Invariant Sections of all of the original documents, unmodified, and list them all as Invariant Sections of your combined work in its license notice, and that you preserve all their Warranty Disclaimers.

The combined work need only contain one copy of this License, and multiple identical Invariant Sections may be replaced with a single copy. If there are multiple Invariant Sections with the same name but different contents, make the title of each such section unique by adding at the end of it, in parentheses, the name of the original author or publisher of that section if known, or else a unique number. Make the same adjustment to the section titles in the list of Invariant Sections in the license notice of the combined work. 2105

In the combination, you must combine any sections Entitled “History” in the various original documents, forming one section Entitled “History”; likewise combine any sections Entitled “Acknowledgements”, and any sections Entitled “Dedications”. You must delete all sections Entitled “Endorsements”. 2110

6. COLLECTIONS OF DOCUMENTS

You may make a collection consisting of the Document and other documents released under this License, and replace the individual copies of this License in the various documents with a single copy that is included in the collection, provided that you follow the rules of this License for verbatim copying of each of the documents in all other respects. 2115

You may extract a single document from such a collection, and distribute it individually under this License, provided you insert a copy of this License into the extracted document, and follow this License in all other respects regarding verbatim copying of that document.

7. AGGREGATION WITH INDEPENDENT WORKS

A compilation of the Document or its derivatives with other separate and independent documents or works, in or on a volume of a storage or distribution medium, is called an “aggregate” if the copyright resulting from the compilation is not used to limit the legal rights of the compilation’s users beyond what the individual works permit. When the Document is included in an aggregate, this License does not apply to the other works in the aggregate which are not themselves derivative works of the Document. 2120

If the Cover Text requirement of section 3 is applicable to these copies of the Document, then if the Document is less than one half of the entire aggregate, the Document’s Cover Texts may be placed on covers that bracket the Document within the aggregate, or the electronic equivalent of covers if the Document is in electronic form. Otherwise they must appear on printed covers that bracket the whole aggregate. 2130

8. TRANSLATION

Translation is considered a kind of modification, so you may distribute translations of the Document under the terms of section 4. Replacing Invariant Sections with translations requires special permission from their copyright holders, but you may include translations of some or all Invariant Sections in addition to the original versions of these Invariant Sections. You may include a translation of this License, and all the license notices in the Document, and any Warranty Disclaimers, provided that you also include the original English version of this License and the original versions of those notices and disclaimers. In case of a disagreement between the translation and the original version of this License or a notice or disclaimer, the original version will prevail. 2135

If a section in the Document is Entitled “Acknowledgements”, “Dedications”, or “History”, the requirement (section 4) to Preserve its Title (section 1) will typically require changing the actual title. 2140

9. TERMINATION

You may not copy, modify, sublicense, or distribute the Document except as expressly provided for under this License. Any other attempt to copy, modify, sublicense or distribute the Document is void, and will automatically terminate your rights under this License. However, parties who have received copies, or rights, from you under this License will not have their licenses terminated so long as such parties remain in full compliance.

10. FUTURE REVISIONS OF THIS LICENSE

The Free Software Foundation may publish new, revised versions of the GNU Free Documentation License from time to time. Such new versions will be similar in spirit to the present version, but may differ in detail to address new problems or concerns. See <http://www.gnu.org/copyleft/>.

Each version of the License is given a distinguishing version number. If the Document specifies that a particular numbered version of this License “or any later version” applies to it, you have the option of following the terms and conditions either of that specified version or of any later version that has been published (not as a draft) by the Free Software Foundation. If the Document does not specify a version number of this License, you may choose any version ever published (not as a draft) by the Free Software Foundation.

ADDENDUM: How to use this License for your documents

To use this License in a document you have written, include a copy of the License in the document and put the following copyright and license notices just after the title page:

```
Copyright (c) YEAR YOUR NAME.
Permission is granted to copy, distribute and/or modify this document
under the terms of the GNU Free Documentation License, Version 1.2
or any later version published by the Free Software Foundation;
with no Invariant Sections, no Front-Cover Texts, and no Back-Cover Texts.
A copy of the license is included in the section entitled ‘‘GNU
Free Documentation License’’.
```

If you have Invariant Sections, Front-Cover Texts and Back-Cover Texts, replace the “with ... Texts.” line with this:

```
with the Invariant Sections being LIST THEIR TITLES, with the
Front-Cover Texts being LIST, and with the Back-Cover Texts being LIST.
```

If you have Invariant Sections without Cover Texts, or some other combination of the three, merge those two alternatives to suit the situation.

If your document contains nontrivial examples of program code, we recommend releasing these examples in parallel under your choice of free software license, such as the GNU General Public License, to permit their use in free software.

Glossary of Terms and Acronyms

Glossary

covert channel Covert channel is a type of computer security flaw that allows to transfer information objects between processes that are not supposed to be allowed to communicate by the computer security policy. 2180

customisation file A user defined file which is loaded in the memory space of XtratuM or the partitions. It is used to pass runtime configuration data to the partitions. For example the configuration vector to XtratuM; or runtime parameters to a partition.

error An error is the part of the system state that may cause a subsequent failure: a failure occurs when an error reaches the service interface and alters the service. 2185

failure A failure is an event that occurs when the delivered service deviates from correct service.

fault A fault is the adjudged or hypothesized cause of an error.

hypercall The service (system call) provided by the hypervisor. The services provided are known as para-virtual services. 2190

hypervisor The layer of software that, using the native hardware resources, provides one or more virtual machines (partitions).

i/o port Or peripheral port, is a low level processor address connected to an external peripheral. Some processors map the I/O ports in a designated memory addresses, and is accessed as if it were RAM memory; while others use a special I/O space which requires special processor instructions. 2195

native hardware The existing hardware: processor, interrupt, clock, etc.

para-virtual A virtual object that resembles, but with a different interface, the native object.

partition Also known as “virtual machine” or “domain”. It refers to the environment created by the hypervisor to execute user code.

partition code Also known as “guest”. Is the code executed inside a partition. Usually, the code is composed of an operating system and a set of processes or threads. Since application code relies on the services provided by the OS, we will assume that the partition code is an operating system (or a real-time operating system). 2200

resident software The booting software that is executed directly in ROM memory right after a system reboot, also referred as boot-loader or firmware. Among other tasks, it is in charge of loading in RAM memory XtratuM and the initial partitions. 2205

side channel A side-channel is any observable information emitted as a byproduct of the physical implementation or operation of the system.

2210

spare time Processor time reserved for future utilisation. Note that idle time is the remaining processor capacity after the current workload has been fully attended.

system partition A partition that has extra capabilities to manage and control the system, and other partitions. Originally these partitions were named “supervisor partitions” but to avoid confusion with the processor modes it was renamed as “system partitions”.

Abbreviated terms

Term	Description
ABI	Application Binary Interface.
APEX	APplication EXecutive.
API	APplication PRogramming Interface.
ARINC	Aeronautical Radio, INC. http://www.arinc.com/
BIOS	Basic Input Output Software.
bps	Bits Per Second.
CC	Common Criteria for Information Technology Security Evaluation.
DMA	Direct Memory Access.
ELF	Executable and Linkable Format.
ESD	Effective Slot Duration.
FIFO	First In First Out.
GPOS	General Purpose Operating System.
HM	Health Monitor.
IMA	Integrated Modular Avionics.
IPC	Inter Partition Communication.
MAF	Major Frame. See cyclic scheduling.
MMU	Memory Management Unit.
PCT	Partition Control Table.
PIT	Partition Information Table.
RSW	Resident SoftWare.
RTEMS	Real-Time Executive for Multiprocessor Systems.
SD	Slot Duration.
ST	Security Target.
TBR	Trap Base Register. A special LEON2 register.
TSC	TSF Scope of Control.
TSO	Total Storage Ordering.
UART	Universal Asynchronous Receiver Transmitter. A serial port.
VMM	Virtual Machine Monitor (hypervisor).
WCET	Worst Case Execution Time.
WIM	Window Invalid Mask. A special LEON2 register.
XAL	XtratuM Abstraction Layer.
XEF	XtratuM Executable Format.
XM_CF	XML XtratuM configuration file. It can also be named as XM_CF.xml to remind that it is an XML file.
XM_CT.bin	The compiled binary version of the XM_CF configuration file.
XML	eXtended Markup Language.

Index

- 2215 bare-application, 39
- boot, 3, 6–8, 48, 75, 76
- channel, 13, 27, 47, 77, 85
- communication port, 8, 13, 18, 22, 27, 89
- component, 28, 71, 76
- 2220 configuration file, 8, 10, 13, 18, 22, 47, 49, 56, 63, 83, 91
- container, 25, 71–73, 75, 76
- context switch, 11, 12, 53
- customisation, 23
- 2225 extended
 - virtual, 14
- health monitor, 6, 14, 27, 49, 51, 85
- hypercall, 39
- I/O server, 18
- 2230 initialise, 39
- integrator, 23
- interrupt, 8, 18
 - hardware, 7, 19
 - native, 22
 - 2235 timer, 11
 - virtual, 22
- major time frame, 9
- menuconfig, 26, 79
- message, 13
- 2240 mode change, 12
- para-virtual, 39
- partition
 - normal, 8
 - standard, 22
 - 2245 system, 6, 8, 13, 21, 22
- PCT, 14, 67
- plan
 - initial, 12
 - maintenance, 12
- 2250 port
 - I/O, 27
 - queuing, 13, 18, 22, 47
 - sampling, 13, 18, 22, 47
- reset, 12, 46
- cold, 7 2255
- hardware, 7
- warm, 7
- resident software, 6, 23, 25, 75, 85
- rswbuild, 96
 - DESCRIPTION, 96 2260
 - SYNOPSIS, 96
 - USAGE EXAMPLES, 96
- scheduling, 3, 6, 8–11, 22, 27, 77
 - cyclic, 3, 9, 10, 22
- state 2265
 - partition, 7
 - system, 6
- time slot, 8
- time slot, 9
- trap table, 18, 19 2270
 - virtual, 19, 49
- xmcparser, 92
 - DESCRIPTION, 92
 - OPTIONS, 92
 - SYNOPSIS, 92 2275
- xmeformat, 92
 - DESCRIPTION, 93
 - SYNOPSIS, 92
 - USAGE EXAMPLES, 93
- xmpack, 94 2280
 - DESCRIPTION, 95
 - SYNOPSIS, 95
 - USAGE EXAMPLES, 96