

Virtualization for Embedded Systems

Is an open source solution right for you?

11/23/2012

Julia Keffer

Table of Contents ✓

| | |
|--|----------|
| Introduction | 1 |
| What is Virtualization? | 1 |
| Use Cases..... | 2 |
| Operating Systems with Different Run-time Requirements | 2 |
| Isolate Security Conscious Applications | 2 |
| Open Source Compliance | 2 |
| Virtualization Implementations | 2 |
| Key Criteria for a Virtualization Solution | 5 |
| Hardware Support..... | 5 |
| Operating System Support | 5 |
| Resource Allocation and Sharing | 5 |
| Memory Isolation | 6 |
| Processor Scheduling | 6 |
| Guest Communication..... | 6 |
| Size of the Code Base | 6 |
| Open Source Solutions..... | 7 |
| Xen..... | 7 |
| Hardware Support | 8 |
| Operating System Support | 8 |
| Resource Allocation and Sharing..... | 8 |
| Memory Isolation | 9 |

| | |
|--------------------------------------|-----------|
| Processor Scheduling..... | 9 |
| Guest Communication..... | 9 |
| Trusted Computing Base | 10 |
| Xtratum | 10 |
| Hardware Support | 10 |
| Operating System Support | 11 |
| Resource Allocation and Sharing..... | 11 |
| Memory Isolation | 11 |
| Processor Scheduling..... | 11 |
| Guest Communication..... | 11 |
| Trusted Computing Base | 11 |
| OKL4 | 11 |
| Hardware Support | 12 |
| Operating System Support | 12 |
| Resource Allocation and Sharing..... | 12 |
| Memory Isolation | 13 |
| Processor Scheduling..... | 13 |
| Guest Communication..... | 13 |
| Trusted Computing Base | 13 |
| Conclusions | 14 |
| Works Cited | 15 |
| Glossary..... | 16 |

Table of Figures ✓

| | |
|--|----|
| Figure 1 – Non –Virtualized and Virtualized Computer (1) | 1 |
| Figure 2- Paravirtualized System (1)..... | 3 |
| Figure 3- Xen Hypervisor (5) | 8 |
| Figure 4 - PCI Pass-Through (5)..... | 9 |
| Figure 5 – Xtratum Architecture (11) | 10 |
| Figure 6 - OKL4 Architecture | 12 |
| Figure 7- OKL4 IPC Model (14)..... | 13 |

Introduction

Embedded computers are part of our everyday lives, from smart phones, to cars, to gaming consoles. Virtualization was predominantly used first in the server market, but today it has come to the embedded computer. Undoubtedly, you have used an embedded computer that employs virtualization technology.

This paper explains what virtualization is, how it is implemented, and how it is applied in embedded applications. It examines a set of criteria for choosing a virtualization solution and evaluates three open source implementations against each of the criteria.

What is Virtualization?

The Computer Desktop Encyclopedia defines a virtual machine as “An operating system that runs like a “machine within a machine”, and functions as if it owned the entire computer, is referred to as a virtual machine (VM). The operating systems in each VM partition are called guest operating systems or partitions, and they communicate with the hardware via the virtual machine control program called a virtual machine monitor (VMM), which is also referred to as a hypervisor. It “virtualizes” the hardware for each guest operating system”.(1)

This paper will use the terms hypervisor and guest to refer to the VMM and the guest operating system, respectively.

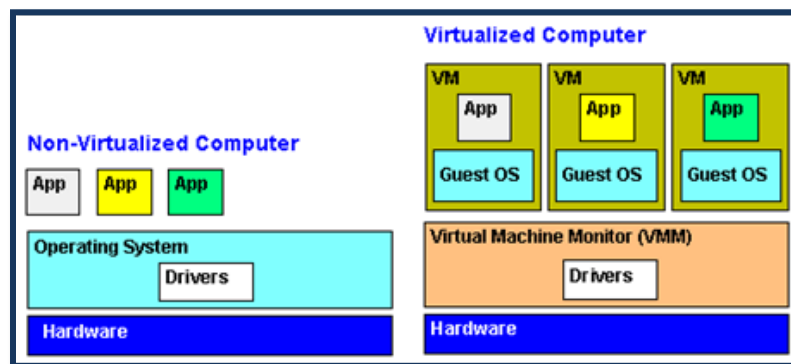


Figure 1 – Non –Virtualized and Virtualized Computer (1)

The hypervisor operates in a privileged environment, also referred to as kernel mode or supervisor mode, where it has access to low level system calls to mediate resource access.

The entire guest operating system operates in a non-privileged environment, also referred to as user mode. Guests communicate with the hardware only through the hypervisor and cannot use low level system calls directly.

The number of guest ~~operating systems~~ that can run on a single hardware platform is constrained by the available hardware resources, typically the amount of memory. Each guest can run a different operating system.

Use Cases

Virtualization is useful either to consolidate multiple computer systems on the same hardware to reduce costs, or to isolate programs running on the same hardware. This section describes three cases where it is useful to run multiple isolated operating systems on the same hardware in embedded systems.

Operating Systems with Different Run-time Requirements

Virtualization provides the ability to run different types of operating systems on the same hardware, such as a full featured OS for user interface functions, and a real-time OS for time critical applications. In an automobile, the computer that controls the anti-lock brake system has real-time requirements, while the infotainment system does not. Previously, an automobile used two different computers, whereas with virtualization, both systems can run on the same hardware, reducing costs.

Isolate Security Conscious Applications

Virtualization can isolate security-conscious applications from insecure applications. In a smart phone, if an application introduces a computer virus, it is necessary to protect the environment where the wireless protocol stack resides to ensure that the system can still make phone calls. One way to do this is to run each of these components inside separate operating systems in a virtual environment.

Open Source Compliance

Open source licenses typically allow proprietary code to interact with open source code if the two communicate only via a messaging interface. If the two types of code run in separate operating systems, the hypervisor fulfills this requirement.

Virtualization Implementations

There are different ~~ways to implement a virtual system~~. The models discussed are type 1, or bare metal hypervisors, where the hypervisor software runs between the hardware and the guest ~~operation system~~. We will discuss three ~~variants~~: full virtualization, paravirtualization, and a microkernel.

With full virtualization, a whole system is emulated (BIOS, disk, processor, network interface, etc.) and a guest ~~operating system~~ runs unmodified on a hypervisor that provides the abstraction of the underlying computer system. The guest ~~operating system~~ is not aware of the hypervisor.

In this model, the hypervisor intercepts hardware access instructions from the guest ~~operating systems~~ and invokes the instructions on behalf of the guest. Full virtualization requires hardware extensions in the computer processor, such as Intel's VT-x technology.

In a paravirtualized system, the guest ~~operation system~~ requires modifications to work in a virtual machine (2) and communicate with the hypervisor. Specifically, some or all of the device drivers in the guest ~~operating system~~ are modified to replace the privileged instructions with direct requests to the hypervisor, which are referred to as hypercalls. ✓

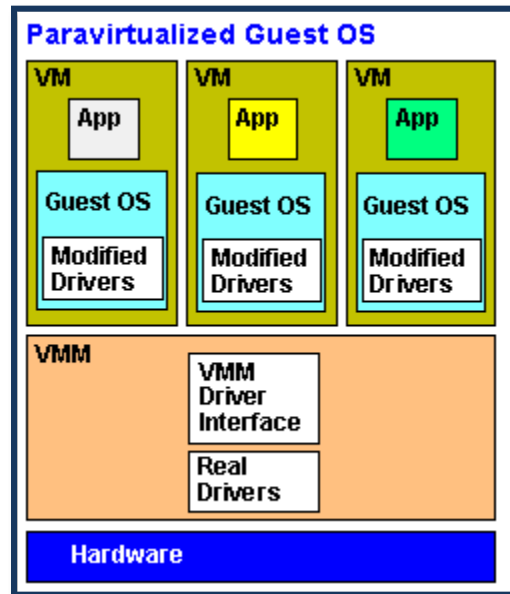



Figure 2- Paravirtualized System (1)

In theory, any operating system which provides access to the source code, such as Linux, can be paravirtualized. Proprietary operating systems, such as Windows, cannot be paravirtualized. Paravirtualization is typically used in systems with hardware that does not support virtualization, although it may still be advantageous to use it for performance reasons. Full virtualization is often not as efficient as paravirtualization, because of the extra step required to intercept the privileged instructions.

Although not originally designed for virtualization, it is possible to implement virtualization on top of a microkernel. A microkernel is a reduced version of a regular operating system kernel which provides a set of policies and mechanisms to access hardware resources.

Any component can run on top of the microkernel. An operating system is a type of component, but it can run alongside a standalone application, for example, a special device driver, that runs directly on top of the microkernel.

As with a hypervisor, the microkernel software runs in kernel mode, between the hardware and the guest ~~operating system~~. Unlike a hypervisor, a microkernel does not perform the instructions on behalf of the guest.

Instead of calling the privileged instructions directly,  it forwards the request to a user mode virtualization component, which interprets the request. The component may reside inside a guest or it may be a standalone component.

The mechanism used to forward the request is called inter-process communication (IPC). As with paravirtualization, the device drivers in the guest ~~operating system~~ must replace hardware access instructions with IPC messages.

Key Criteria for a Virtualization Solution

Before considering virtualization in an embedded system, there are a number of factors you should consider. These include:

- hardware support
- operating system support
- resource allocation and sharing
- memory isolation
- processor scheduling
- guest communication
- size of the code base

Hardware Support


The hypervisor or microkernel software runs directly on top of the hardware, and therefore must support the instructions required by the hardware architecture. Common architectures in embedded systems are x86, ARM, PowerPC, and Sparc. The x86 processor is typically used in industrial and medical applications. Smart phones and tablets almost exclusively use ARM processors. The Sparc architecture is common in military and avionics systems. Gaming consoles use PowerPC architectures.

To support full virtualization, the hardware must include virtualization extensions. Intel and AMD both include virtualization support in their x86 processors. The ARM Cortex A15 and A7 processors also support virtualization.

Full memory isolation requires a memory management unit (MMU).

Operating System Support

Hardware and operating system support are closely related. As mentioned previously, some operating systems, such as Microsoft Windows, cannot be paravirtualized, and therefore need hardware support.

If paravirtualization is necessary or desirable, access to the operating system code is required, which is typically available in open source operating systems, such as all variants of Linux, FreeBSD, NetBSD, and OpenSolaris.  Some Linux operating systems already include the paravirtualized drivers (3).

Resource Allocation and Sharing

Guests must share some hardware resources, such as disks and network interfaces. An application in the system may require dedicated access to a particular device, such as a USB port, which means that the hypervisor must provide a mechanism to assign exclusive access to the device to a specific guest ~~operating system~~.

Memory Isolation

The memory allocation scheme must ensure that guests cannot access memory outside their own address range. It is important to note that any truly secure implementation requires hardware support by a memory management unit (MMU) to guard against a malicious device driver that uses direct memory access (DMA). Both Intel and AMD processors have MMU support, as do some ARM processors.


Processor Scheduling

Execution isolation is important to prevent a rogue application on a guest from monopolizing the CPU, which essentially functions as a denial-of-service attack on the rest of the guests. If one of the guests requires a real-time response, the hypervisor must use a scheduling algorithm that can assign it a higher priority. It is also desirable to have a way of ensuring that lower priority tasks in one guest do not preempt higher priority tasks in another guest.

Guest Communication

Guests may want to exchange information. For example, a component may need to provide status information for a user interface to display. If there is a mechanism to enable guests to communicate, the solution must prevent a security breach using this mechanism. Any risk typically results from the mechanism the guests use to store the data to exchange.

Size of the Code Base

The size of the code base that implements virtualization affects system robustness. All code has a certain number of defects and the smaller the code base, the fewer defects there are likely to be. Because the code runs in privileged mode, it must be possible to contain the faults to the virtualization code without affecting the guest  operating systems. The code with access to the privileged instructions is referred to as the trusted computing base (TCB).

Open Source Solutions ✓

Open source solutions are appealing for a number of reasons, but mainly because the code is free and there is the flexibility to change it according to the needs of the system. An ideal open source solution has active development and community members willing to informally support users.

One potential drawback of open source is the requirements under GNU GPL¹ to release the source code for any derived work. If it is necessary to make proprietary changes to the virtualization code, open source code may not be an appropriate solution.

The open source solutions described in this section are licensed under either GPL or a proprietary license with the same conditions as GPL. Each implements virtualization in a slightly different way, each with advantages and disadvantages.

Xen

Xen is a bare metal hypervisor. There are two parts to the hypervisor implementation. The hypervisor code that runs directly on top of the hardware is responsible for virtualizing the CPU, memory, and input/output (I/O) control, including interrupt handling. ✓

A special paravirtualized guest (referred to as Domain 0) has privileged access to the hardware. It manages processor and memory sharing, network and disk access, and communication between guests. Domain 0 can run any paravirtualized operating system, but it is typically a variant of Linux, as many Linux distributions include native Xen support.

Xen supports both full virtualization and paravirtualization. Xen has an active contribution community. Refer to (4) for more information.

¹ Refer to (17) for more information.

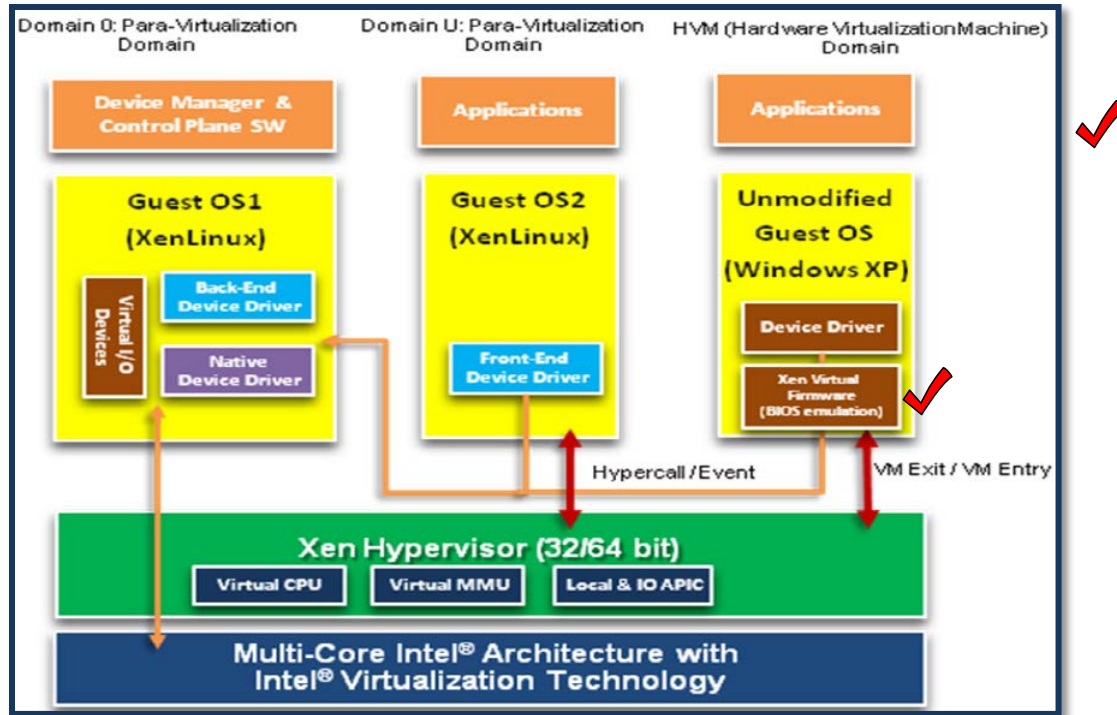


Figure 3- Xen Hypervisor (5)

Hardware Support - Xen can run on x86 processors from Intel and AMD. Xen support for ARM processors is a project led by Samsung which delivers and maintains Xen support for a range of ARM processors (ARM v5 - v7) for mobile devices. The project is also working on problems such as solving real-time guarantees in a virtualized environment and multi-processor support. Refer to (6) for information. An experimental version of Xen which uses the virtualization support introduced for the ARM Cortex A15 is underway. An experimental project to port Xen to PowerPC was abandoned.

Operating System Support – Xen supports any guest ~~operating system~~ that can be paravirtualized. It supports full virtualization for any guest ~~operating system~~ that runs on Intel or AMD x86 hardware with virtualization extensions. Many Linux distributions include the virtual device drivers to support paravirtualization on Xen.

Resource Allocation and Sharing – Domain 0 mediates access to I/O devices by receiving requests from the guests on a virtual channel. Xen also supports a PCI pass-through interface to allow guests direct and exclusive hardware access to PCI devices, such as the network interface. This hardware access method requires Intel VT-x or AMD-V hardware support, and can be used with paravirtualized and fully virtualized guests. The guest must have a native device driver for the device and the Domain 0 guest must have a “pciback” version of the driver. ✓

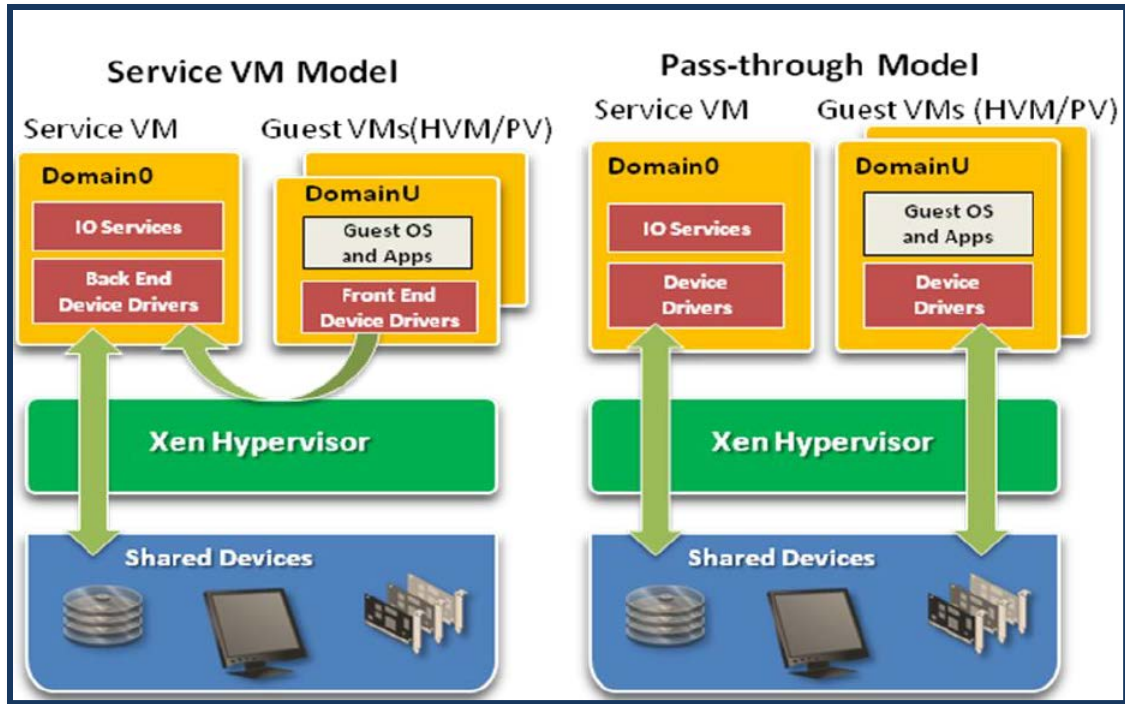


Figure 4 - PCI Pass-Through (5)

The PCI pass-through feature must be enabled in both the BIOS and in Xen and is a potential security risk if a guest runs a malicious application or has defects in its device driver.

Additional pass-through support is available for USB devices and selected graphics devices. A mechanism called SR-IOV allows devices to be assigned to guests but shared among them. Refer to (5) and (7) for information about pass-through functionality.

Memory Isolation – Xen assigns a static area of memory to each guest. The hypervisor uses shadow pages to translate virtual memory access requests from the guests into the physical address. It tracks which guest owns the memory to enforce isolation (8). Newer releases of Xen implement a feature to allow identical guest operating systems to share physical memory for common binaries and libraries. This feature is still in the beta stage and does not have security support. Xen uses shared memory to implement guest communication. Domain 0 manages grant tables which grant access to guests on a per page basis to ensure safe memory sharing (9).

Processor Scheduling – Xen has settings to configure the CPU usage across guests. It load balances across CPUs using a weight and cap (limit) for each guest. It is possible to configure a guest to use only to a specified set of CPUs; however, a guest cannot be assigned exclusively to a single CPU. CPU configuration takes into account physical devices and hyperthreading.

Guest Communication – Guests communicate through a virtual network interface in domain 0, which implements routing and bridging functionality. Communication is based on standard networking protocols. The

default implementation of guest communication requires significant overhead. The XenLoop and XenSockets projects attempt to address this issue (10).

Trusted Computing Base – Although the portion of Xen that resides on top of the hardware is small, domain 0 is a complete operating system. The TCB for Xen is quite large and susceptible to defects in many areas. A newer version of Xen on ARM is experimenting with moving device drivers outside of domain 0 to isolate them from the rest of the computing base.

Xtratum

Xtratum is a bare metal hypervisor that runs in privileged mode and virtualizes the CPU, memory, interrupts, and any devices that endanger isolation. A guest ~~operating system~~, which Xtratum documents refer to as a partition, must be paravirtualized to run on Xtratum and replace calls to privileged instructions with hypercalls.

Xtratum supports a special system guest which can have extra rights to manage and monitor system resources, such as stopping, starting or resetting partitions, using a special set of hypercalls. The access rights of a system guest are set in the Xtratum configuration file. System guests do not have direct hardware access.

The scheduling and IPC mechanisms are modeled on the ARINC (Avionics Application Standard Software Interface) 653 standards, although its goal does not include compliance to the specification. The Universidad Politecnica de Valencia in Spain developed Xtratum. Refer to (11) for more information. ✓

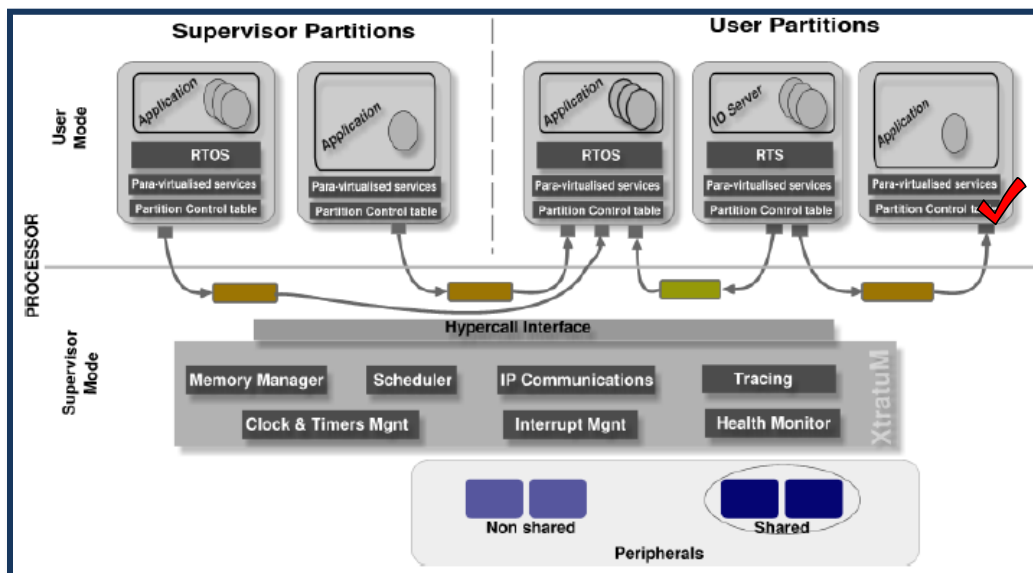


Figure 5 – Xtratum Architecture (11)

Hardware Support – Xtratum supports the LEON3 (Sparc V8) and Intel Itanium-64 processor, neither of which have virtualization extensions. A theoretical paper was published about porting Xtratum to PowerPC but was never implemented (12).

Operating System Support – Any paravirtualized guest ~~operating system~~ can run on Xtratum.

Resource Allocation and Sharing – I/O ports and interrupts that the hypervisor does not manage are assigned exclusively to guest ~~operating systems~~ in the configuration file. The device driver resides in the guest. To share devices, the system designer must implement an I/O server partition which receives requests from other guests via IPC and processes them according to its policy configuration.

Memory Isolation – The configuration file defines the memory area statically assigned to the guest. There are no shared memory regions. Xtratum will run without an MMU, in which case, there is a risk of unauthorized memory access.

Processor Scheduling – Xtratum uses a fixed, cyclical scheduling algorithm based on the timeslot and duration settings for each guest in the configuration file. Each guest may define multiple scheduling plans and can notify the hypervisor to switch plans using a hypercall. For example, a guest may want to switch into a maintenance mode when it has only low priority tasks to do, freeing the processor for use by other guests with higher priority tasks. A system guest can also change the scheduling plan of a normal guest.

Guest Communication – The hypervisor implements a port-based communication mechanism. Guests send and receive messages from each other or the hypervisor on a channel that links two ports. The protocol is specific to the sending and receiving parties. Both broadcast and direct messaging modes are available. Channels, ports, maximum message sizes, and maximum number of messages (queuing ports) are defined in the configuration file. Data exchange relies on buffer copying mechanisms, as there are no shared memory regions.

Trusted Computing Base – The critical code for Xtratum is limited to the small hypervisor code base. It uses a health monitor feature to detect and react to errors to contain them within the proper scope: process, guest, hypervisor, or firmware.

OKL4

OKL4 3.0² is a microkernel implementation of virtualization. The microkernel runs on top of the hardware in kernel mode and uses inter-process messaging (IPC) to mediate requests for interrupts and device drivers between guest ~~operating systems~~. The microkernel does not provide system services; these are implemented by separate components running in user mode.

A separate resource and policy model, which runs outside the microkernel in user space, holds the configuration of the CPU scheduling policy and the memory allocation. The microkernel runs in privileged mode and the guests run in user mode. Guests must replace hardware access instructions in their device drivers with IPC messages. Open Kernel Labs, which is owned by General Dynamics, sponsors the OKL4 project. Refer to (13) for more information.

² Not to be confused with the OKL4 4.0 microvisor, which requires a commercial license.

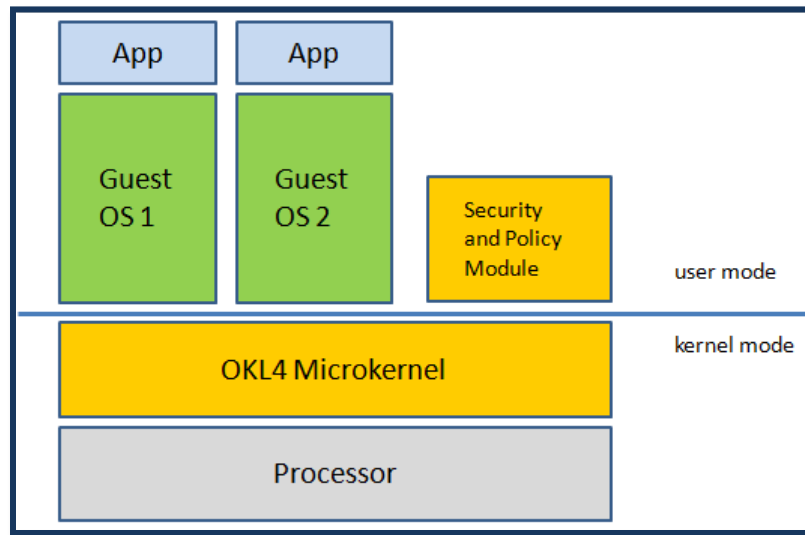


Figure 6 - OKL4 Architecture

Hardware Support - OKL4 supports ARMv5/v6 and Intel i386 processors, none of which have virtualization extensions. OKL4 must run on a processor with a memory management unit (MMU).

Operating System Support – Any paravirtualized guest ~~operating system~~ can run on OKL4. Open Kernel Labs provides a paravirtualized version of Linux to use as a guest ~~operating system~~.

Resource Allocation and Sharing – OKL4 mediates device access using IPC messages. It relays requests for device access from a guest's virtual driver to the physical driver, which may be either a standalone driver, or reside in another guest ~~operating system~~. The policy module controls which guest can drive a particular device by mapping device registers.

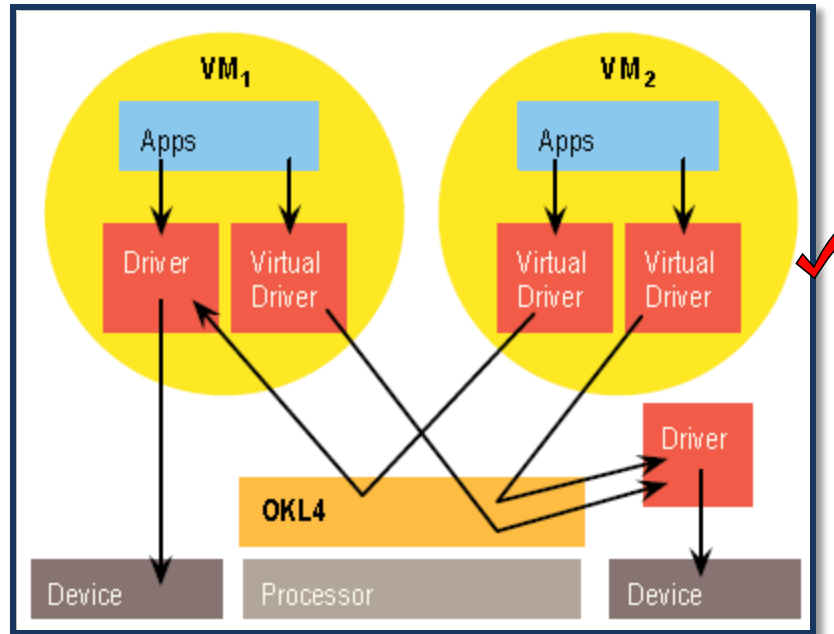


Figure 7- OKL4 IPC Model (14)

Memory Isolation – The memory allocated for each guest is statically configured in the resource and policy module. There are also configuration settings for shared memory regions and policy settings to determine which guests can access the shared regions. The policy module has a monopoly over operations that consume kernel memory; it can control which guest is allowed to consume such kernel resources to guard against denial-of-service attacks on the system, for example, by a rogue guest kernel.

Processor Scheduling – The system designer configures the CPU scheduling algorithm in the security and policy module. The policy can schedule the entire the guest ~~operating system~~ an individual process within a ✓ guest. A global policy which assigns process priorities across the guests ensures that lower priority maintenance tasks in one guest do not block higher priority processes in another guest. The range of priorities the guest can use is restricted to ensure that it does not monopolize the system.

Guest Communication – Guest communication is the heart of the OKL4 implementation and performance has been highly optimized. In addition to device access requests to the microkernel, guest ✓ can exchange information based on an agreed-upon protocol. They ✓ may make use of the shared memory regions to facilitate buffer allocation and access to shared data.

Trusted Computing Base – The microkernel code base is very small and therefore, less likely to have defects. Much effort has gone toward providing a formal proof of correctness of the kernel. Its separation from the policy module, which is located in user space, allows flexibility while limiting the code that has access to privileged activities.

Conclusions ✓

The main consideration for choosing a virtualization solution for an embedded system is the hardware platform. As we can see, each virtualization solution supports a different range of hardware vendors. With an open source solution, while it is possible to implement one of the solutions on a different platform, it may require significant effort.

Closely tied to the hardware solution is the choice of operating system. Most Linux distributions work well on any of the solutions, although Linux does not have real-time support. ✓ To save time and effort, you may also want to choose an operating system that has paravirtualized drivers already available.

Once you choose the hardware and operating system, the needs of the application drive the choice of the virtualization solution. You need to consider how the implementation of the virtualization solution affects the design of your application. You may need to implement an I/O server to share devices, or you may be able to take advantage of hardware support to have exclusive access to a device. A particular CPU sharing mechanism may be a better fit, depending on the scheduling constraints of your application. The use of a shared memory region to exchange data between guests may be a benefit or a concern.

Choosing the right solution is important to ensure that you can create an application that is easy to design and maintain.



Works Cited

1. Computer Desktop Encyclopedia. [Online] <http://www.computerlanguage.com/>.
2. Paravirtualization explained. *Search Server Virtualization*. [Online] TechTarget, 2007. <http://searchservirtualization.techtarget.com/tip/Paravirtualization-explained>.
3. Paravirtualization. *Wikipedia*. [Online] <http://en.wikipedia.org/wiki/Paravirtualization>.
4. Xen Wiki. [Online] <http://wiki.xen.org>.
5. **Amit Aneja**. *Designing Embedded Virtualization Intel Platform*. s.l. : Intel Corporation, 2011.
6. Xen ARM Wiki. *Xen ARM Wiki*. [Online] Samsung Corp. <http://wiki.xen.org/wiki/XenARM>.
7. **Intel Corporation**. *Intel® Virtualization Technology for Directed I/O Architecture Description*. s.l. : Intel Corporation, 2011.
8. Memory Allocation - Xen 3.0 Virtualization Interface Guide. *Linuxtopia*. [Online] http://www.linuxtopia.org/online_books/linux_virtualization/xen_3.0_interface_guide/linux_virtualization_xen_interface_9.html.
9. Inter-Domain Communication - Xen 3.0 Virtualization Interface Guide. *Linuxtopia*. [Online] http://www.linuxtopia.org/online_books/linux_virtualization/xen_3.0_interface_guide/linux_virtualization_xen_interface_52.html.
10. **Tomlinson, Allan and Gebhardt, Carl**. *Challenges for Inter Virtual Machine Communication*. Mathematics, Royal Holloway, University of London. Surrey : s.n., 2010.
11. **Miguel Masmano, Ismael Ripoll, Alfons Crespo**. *XtratuM Hypervisor for LEON3 Volume2: User Manual*. s.l. : Universidad Politecnica de Valencia, February 2011.
12. **Zhou, Rui**. *Partitioned System with XtratuM on PowerPC*. Universi. s.l. : Universidad Politecnica de Valencia, 2009.
13. Open Kernel Labs Community Wiki. [Online] <http://wiki.ok-labs.com/>.
14. **Gernot Heiser**. *Virtualization for Embedded Systems*. s.l. : Open Kernel Labs, 2007.
15. Xen Web site. [Online] Citrix. <http://wiki.xen.org>.
16. Xen. *Wikipedia*. [Online] <http://en.wikipedia.org/wiki/Xen>.
17. GNU General Public License. [Online] <http://www.gnu.org/licenses/gpl.html>.

Glossary

AMD – Advanced Micro Devices

ARINC – Avionics Application Standard Software Interface

ARM – Originally Acorn RISC (Reduced Instruction Set Computer) Machine

BIOS – Basic Input/Output System

CPU – Central Processing Unit

DMA – Directory Memory Access

GNU – Gnu's Not Unix

GPL – General Public License

I/O – Input/Output

IPC – Inter-Process Communication

MMU – Memory Management Unit

PCI – Peripheral Component Interface

TCB – Trusted Computing Base

USB – Universal Serial Bus

VM – Virtual Machine

VMM – Virtual Machine Monitory