



Users' Manual

Xen v3.0

DISCLAIMER: This documentation is always under active development and as such there may be mistakes and omissions — watch out for these and please report any you find to the developers' mailing list, xen-devel@lists.xensource.com. The latest version is always available on-line. Contributions of material, suggestions and corrections are welcome.

Xen is Copyright ©2002-2005, University of Cambridge, UK, XenSource Inc., IBM Corp., Hewlett-Packard Co., Intel Corp., AMD Inc., and others. All rights reserved.

Xen is an open-source project. Most portions of Xen are licensed for copying under the terms of the GNU General Public License, version 2. Other portions are licensed under the terms of the GNU Lesser General Public License, the Zope Public License 2.0, or under “BSD-style” licenses. Please refer to the COPYING file for details.

Xen includes software by Christopher Clark. This software is covered by the following licence:

Copyright (c) 2002, Christopher Clark. All rights reserved.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

- Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.
- Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.
- Neither the name of the original author; nor the names of any contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

Contents

1	Introduction	1
1.1	Usage Scenarios	1
1.2	Operating System Support	2
1.3	Hardware Support	2
1.4	Structure of a Xen-Based System	3
1.5	History	3
1.6	What's New	4
I	Installation	5
2	Basic Installation	7
2.1	Prerequisites	7
2.2	Installing from Binary Tarball	8
2.3	Installing from RPMs	8
2.4	Installing from Source	8
2.4.1	Obtaining the Source	8
2.4.2	Building from Source	9
2.4.3	Custom Kernels	9
2.4.4	Installing Generated Binaries	10
2.5	Configuration	10
2.5.1	GRUB Configuration	10
2.5.2	Serial Console (optional)	11
2.5.3	TLS Libraries	14
2.6	Booting Xen	14
3	Booting a Xen System	15
3.1	Booting Domain0	15
3.2	Booting Guest Domains	16
3.2.1	Creating a Domain Configuration File	16
3.2.2	Booting the Guest Domain	16
3.3	Starting / Stopping Domains Automatically	17

II	Configuration and Management	19
4	Domain Management Tools	21
4.1	Xend	21
4.1.1	Logging	22
4.1.2	Configuring Xend	22
4.2	Xm	23
4.2.1	Basic Management Commands	23
4.2.2	Domain Scheduling Management Commands	24
5	Domain Configuration	25
5.1	Configuration Files	25
5.2	Network Configuration	26
5.2.1	Xen virtual network topology	26
5.2.2	Xen networking scripts	27
5.3	Driver Domain Configuration	27
5.3.1	PCI	27
5.4	Support for virtual Trusted Platform Module (vTPM)	30
6	Storage and File System Management	33
6.1	Exporting Physical Devices as VBDs	33
6.2	Using File-backed VBDs	34
6.2.1	Loopback-mounted file-backed VBDs (deprecated)	35
6.3	Using LVM-backed VBDs	36
6.4	Using NFS Root	37
7	CPU Management	39
8	Migrating Domains	41
8.1	Domain Save and Restore	41
8.2	Migration and Live Migration	41
9	Securing Xen	43
9.1	Xen Security Considerations	43
9.2	Driver Domain Security Considerations	43
9.3	Security Scenarios	45
9.3.1	The Isolated Management Network	45
9.3.2	A Subnet Behind a Firewall	45
9.3.3	Nodes on an Untrusted Subnet	45
10	sHype/Xen Access Control	47
10.1	Overview	48
10.2	Xen Workload Protection Step-by-Step	49

10.2.1	Configuring/Building sHype Support into Xen	49
10.2.2	Creating A WLP Policy in 3 Simple Steps with ezPolicy . . .	50
10.2.3	Deploying a WLP Policy	51
10.2.4	Labeling Domains	52
10.2.5	Labeling Resources	53
10.2.6	Testing The Xen Workload Protection	54
10.3	Xen Access Control Policy	56
10.3.1	Policy Header and Policy Name	56
10.3.2	Simple Type Enforcement Policy Component	58
10.3.3	Chinese Wall Policy Component	58
10.3.4	Security Labels	59
10.3.5	Tools For Creating sHype/Xen Security Policies	62
10.4	Current Limitations	62
10.4.1	Network Traffic	62
10.4.2	Resource Access and Usage Control	63
10.4.3	Domain Migration	63
10.4.4	Covert Channels	63
III	Reference	65
11	Build and Boot Options	67
11.1	Top-level Configuration Options	67
11.2	Xen Build Options	67
11.3	Xen Boot Options	68
11.4	XenLinux Boot Options	70
12	Further Support	71
12.1	Other Documentation	71
12.2	Online References	71
12.3	Mailing Lists	72
A	Unmodified (VMX) guest domains in Xen with Intel®Virtualization Technology (VT)	73
A.1	Building Xen with VT support	73
A.2	Configuration file for unmodified VMX guests	74
A.3	Creating virtual disks from scratch	76
A.3.1	Using physical disks	76
A.3.2	Using disk image files	76
A.3.3	Install Windows into an Image File using a VMX guest	78
A.4	VMX Guests	79
A.4.1	Editing the Xen VMX config file	79

A.4.2	Creating VMX guests	79
A.4.3	Mouse issues, especially under VNC	79
A.4.4	USB Support	82
A.4.5	Destroy VMX guests	84
A.4.6	VMX window (X or VNC) Hot Key	84
A.4.7	Save/Restore and Migration	84
B	Vnets - Domain Virtual Networking	85
B.1	Example	86
B.2	Installing vnet support	87
C	Glossary of Terms	89

Chapter 1

Introduction

Xen is an open-source *para-virtualizing* virtual machine monitor (VMM), or “hypervisor”, for the x86 processor architecture. Xen can securely execute multiple virtual machines on a single physical system with close-to-native performance. Xen facilitates enterprise-grade functionality, including:

- Virtual machines with performance close to native hardware.
- Live migration of running virtual machines between physical hosts.
- Up to 32 virtual CPUs per guest virtual machine, with VCPU hotplug.
- x86/32, x86/32 with PAE, and x86/64 platform support.
- Intel Virtualization Technology (VT-x) for unmodified guest operating systems (including Microsoft Windows).
- Excellent hardware support (supports almost all Linux device drivers).

1.1 Usage Scenarios

Usage scenarios for Xen include:

Server Consolidation. Move multiple servers onto a single physical host with performance and fault isolation provided at the virtual machine boundaries.

Hardware Independence. Allow legacy applications and operating systems to exploit new hardware.

Multiple OS configurations. Run multiple operating systems simultaneously, for development or testing purposes.

Kernel Development. Test and debug kernel modifications in a sand-boxed virtual machine — no need for a separate test machine.

Cluster Computing. Management at VM granularity provides more flexibility than

separately managing each physical host, but better control and isolation than single-system image solutions, particularly by using live migration for load balancing.

Hardware support for custom OSes. Allow development of new OSes while benefiting from the wide-ranging hardware support of existing OSes such as Linux.

1.2 Operating System Support

Para-virtualization permits very high performance virtualization, even on architectures like x86 that are traditionally very hard to virtualize.

This approach requires operating systems to be *ported* to run on Xen. Porting an OS to run on Xen is similar to supporting a new hardware platform, however the process is simplified because the para-virtual machine architecture is very similar to the underlying native hardware. Even though operating system kernels must explicitly support Xen, a key feature is that user space applications and libraries *do not* require modification.

With hardware CPU virtualization as provided by Intel VT and AMD SVM technology, the ability to run an unmodified guest OS kernel is available. No porting of the OS is required, although some additional driver support is necessary within Xen itself. Unlike traditional full virtualization hypervisors, which suffer a tremendous performance overhead, the combination of Xen and VT or Xen and Pacifica technology complement one another to offer superb performance for para-virtualized guest operating systems and full support for unmodified guests running natively on the processor. Full support for VT and Pacifica chipsets will appear in early 2006.

Paravirtualized Xen support is available for increasingly many operating systems: currently, mature Linux support is available and included in the standard distribution. Other OS ports—including NetBSD, FreeBSD and Solaris x86 v10—are nearing completion.

1.3 Hardware Support

Xen currently runs on the x86 architecture, requiring a “P6” or newer processor (e.g. Pentium Pro, Celeron, Pentium II, Pentium III, Pentium IV, Xeon, AMD Athlon, AMD Duron). Multiprocessor machines are supported, and there is support for Hyper-Threading (SMT). In addition, ports to IA64 and Power architectures are in progress.

The default 32-bit Xen supports up to 4GB of memory. However Xen 3.0 adds support for Intel’s Physical Addressing Extensions (PAE), which enable x86/32 machines to address up to 64 GB of physical memory. Xen 3.0 also supports x86/64 platforms such

as Intel EM64T and AMD Opteron which can currently address up to 1TB of physical memory.

Xen offloads most of the hardware support issues to the guest OS running in the *Domain 0* management virtual machine. Xen itself contains only the code required to detect and start secondary processors, set up interrupt routing, and perform PCI bus enumeration. Device drivers run within a privileged guest OS rather than within Xen itself. This approach provides compatibility with the majority of device hardware supported by Linux. The default XenLinux build contains support for most server-class network and disk hardware, but you can add support for other hardware by configuring your XenLinux kernel in the normal way.

1.4 Structure of a Xen-Based System

A Xen system has multiple layers, the lowest and most privileged of which is Xen itself.

Xen may host multiple *guest* operating systems, each of which is executed within a secure virtual machine. In Xen terminology, a *domain*. Domains are scheduled by Xen to make effective use of the available physical CPUs. Each guest OS manages its own applications. This management includes the responsibility of scheduling each application within the time allotted to the VM by Xen.

The first domain, *domain 0*, is created automatically when the system boots and has special management privileges. Domain 0 builds other domains and manages their virtual devices. It also performs administrative tasks such as suspending, resuming and migrating other virtual machines.

Within domain 0, a process called *xend* runs to manage the system. Xend is responsible for managing virtual machines and providing access to their consoles. Commands are issued to xend over an HTTP interface, via a command-line tool.

1.5 History

Xen was originally developed by the Systems Research Group at the University of Cambridge Computer Laboratory as part of the XenoServers project, funded by the UK-EPSC.

XenoServers aim to provide a “public infrastructure for global distributed computing”. Xen plays a key part in that, allowing one to efficiently partition a single machine to enable multiple independent clients to run their operating systems and applications in an environment. This environment provides protection, resource isolation and accounting. The project web page contains further information along with pointers to papers and technical reports: <http://www.cl.cam.ac.uk/xeno>

Xen has grown into a fully-fledged project in its own right, enabling us to investigate interesting research issues regarding the best techniques for virtualizing resources such as the CPU, memory, disk and network. Project contributors now include XenSource, Intel, IBM, HP, AMD, Novell, RedHat.

Xen was first described in a paper presented at SOSP in 2003¹, and the first public release (1.0) was made that October. Since then, Xen has significantly matured and is now used in production scenarios on many sites.

1.6 What's New

Xen 3.0.0 offers:

- Support for up to 32-way SMP guest operating systems
- Intel (Physical Addressing Extensions) PAE to support 32-bit servers with more than 4GB physical memory
- x86/64 support (Intel EM64T, AMD Opteron)
- Intel VT-x support to enable the running of unmodified guest operating systems (Windows XP/2003, Legacy Linux)
- Enhanced control tools
- Improved ACPI support
- AGP/DRM graphics

Xen 3.0 features greatly enhanced hardware support, configuration flexibility, usability and a larger complement of supported operating systems. This latest release takes Xen a step closer to being the definitive open source solution for virtualization.

¹<http://www.cl.cam.ac.uk/netos/papers/2003-xensosp.pdf>

Part I

Installation

Chapter 2

Basic Installation

The Xen distribution includes three main components: Xen itself, ports of Linux and NetBSD to run on Xen, and the userspace tools required to manage a Xen-based system. This chapter describes how to install the Xen 3.0 distribution from source. Alternatively, there may be pre-built packages available as part of your operating system distribution.

2.1 Prerequisites

The following is a full list of prerequisites. Items marked ‘†’ are required by the xend control tools, and hence required if you want to run more than one virtual machine; items marked ‘*’ are only required if you wish to build from source.

- A working Linux distribution using the GRUB bootloader and running on a P6-class or newer CPU.
- † The `iproute2` package.
- † The Linux `bridge-utils`¹ (e.g., `/sbin/brctl`)
- † The Linux `hotplug` system² (e.g., `/sbin/hotplug` and related scripts). On newer distributions, this is included alongside the Linux `udev` system³.
- * Build tools (`gcc v3.2.x` or `v3.3.x`, `binutils`, `GNU make`).
- * Development installation of `zlib` (e.g., `zlib-dev`).
- * Development installation of `Python v2.2` or later (e.g., `python-dev`).
- * `LATEX` and `transfig` are required to build the documentation.

¹Available from <http://bridge.sourceforge.net>

²Available from <http://linux-hotplug.sourceforge.net/>

³See <http://www.kernel.org/pub/linux/utils/kernel/hotplug/udev.html/>

Once you have satisfied these prerequisites, you can now install either a binary or source distribution of Xen.

2.2 Installing from Binary Tarball

Pre-built tarballs are available for download from the XenSource downloads page:

```
http://www.xensource.com/downloads/
```

Once you've downloaded the tarball, simply unpack and install:

```
# tar zxvf xen-3.0-install.tgz
# cd xen-3.0-install
# sh ./install.sh
```

Once you've installed the binaries you need to configure your system as described in Section 2.5.

2.3 Installing from RPMs

Pre-built RPMs are available for download from the XenSource downloads page:

```
http://www.xensource.com/downloads/
```

Once you've downloaded the RPMs, you typically install them via the RPM commands:

```
# rpm -iv rpmname
```

See the instructions and the Release Notes for each RPM set referenced at:

```
http://www.xensource.com/downloads/.
```

2.4 Installing from Source

This section describes how to obtain, build and install Xen from source.

2.4.1 Obtaining the Source

The Xen source tree is available as either a compressed source tarball or as a clone of our master Mercurial repository.

Obtaining the Source Tarball

Stable versions and daily snapshots of the Xen source tree are available from the Xen download page:

```
http://www.xensource.com/downloads/
```

Obtaining the source via Mercurial

The source tree may also be obtained via the public Mercurial repository at:

`http://xenbits.xensource.com`

See the instructions and the Getting Started Guide referenced at:

`http://www.xensource.com/downloads/`

2.4.2 Building from Source

The top-level Xen Makefile includes a target “world” that will do the following:

- Build Xen.
- Build the control tools, including xend.
- Download (if necessary) and unpack the Linux 2.6 source code, and patch it for use with Xen.
- Build a Linux kernel to use in domain 0 and a smaller unprivileged kernel, which can be used for unprivileged virtual machines.

After the build has completed you should have a top-level directory called `dist/` in which all resulting targets will be placed. Of particular interest are the two XenLinux kernel images, one with a “-xen0” extension which contains hardware device drivers and drivers for Xen’s virtual devices, and one with a “-xenU” extension that just contains the virtual ones. These are found in `dist/install/boot/` along with the image for Xen itself and the configuration files used during the build.

To customize the set of kernels built you need to edit the top-level Makefile. Look for the line:

```
KERNELS ?= linux-2.6-xen0 linux-2.6-xenU
```

You can edit this line to include any set of operating system kernels which have configurations in the top-level `buildconfigs/` directory.

2.4.3 Custom Kernels

If you wish to build a customized XenLinux kernel (e.g. to support additional devices or enable distribution-required features), you can use the standard Linux configuration mechanisms, specifying that the architecture being built for is `xen`, e.g:

```
# cd linux-2.6.12-xen0
# make ARCH=xen xconfig
# cd ..
# make
```

You can also copy an existing Linux configuration (`.config`) into e.g. `linux-2.6.12-xen0` and execute:

```
# make ARCH=xen oldconfig
```

You may be prompted with some Xen-specific options. We advise accepting the defaults for these options.

Note that the only difference between the two types of Linux kernels that are built is the configuration file used for each. The “U” suffixed (unprivileged) versions don’t contain any of the physical hardware device drivers, leading to a 30% reduction in size; hence you may prefer these for your non-privileged domains. The “0” suffixed privileged versions can be used to boot the system, as well as in driver domains and unprivileged domains.

2.4.4 Installing Generated Binaries

The files produced by the build process are stored under the `dist/install/` directory. To install them in their default locations, do:

```
# make install
```

Alternatively, users with special installation requirements may wish to install them manually by copying the files to their appropriate destinations.

The `dist/install/boot` directory will also contain the config files used for building the XenLinux kernels, and also versions of Xen and XenLinux kernels that contain debug symbols such as `(xen-syms-3.0.0` and `vmlinux-syms-2.6.12.6-xen0)` which are essential for interpreting crash dumps. Retain these files as the developers may wish to see them if you post on the mailing list.

2.5 Configuration

Once you have built and installed the Xen distribution, it is simple to prepare the machine for booting and running Xen.

2.5.1 GRUB Configuration

An entry should be added to `grub.conf` (often found under `/boot/` or `/boot/grub/`) to allow Xen / XenLinux to boot. This file is sometimes called `menu.lst`, depending on your distribution. The entry should look something like the following:

```
title Xen 3.0 / XenLinux 2.6
  kernel /boot/xen-3.0.gz dom0_mem=262144
  module /boot/vmlinuz-2.6-xen0 root=/dev/sda4 ro console=tty0
```

The kernel line tells GRUB where to find Xen itself and what boot parameters should be passed to it (in this case, setting the domain 0 memory allocation in kilobytes and

the settings for the serial port). For more details on the various Xen boot parameters see Section 11.3.

The module line of the configuration describes the location of the XenLinux kernel that Xen should start and the parameters that should be passed to it. These are standard Linux parameters, identifying the root device and specifying it be initially mounted read only and instructing that console output be sent to the screen. Some distributions such as SuSE do not require the `ro` parameter.

To use an `initrd`, add another module line to the configuration, like:

```
module /boot/my_initrd.gz
```

When installing a new kernel, it is recommended that you do not delete existing menu options from `menu.lst`, as you may wish to boot your old Linux kernel in future, particularly if you have problems.

2.5.2 Serial Console (optional)

Serial console access allows you to manage, monitor, and interact with your system over a serial console. This can allow access from another nearby system via a null-modem (“LapLink”) cable or remotely via a serial concentrator.

Your system’s BIOS, bootloader (GRUB), Xen, Linux, and login access must each be individually configured for serial console access. It is *not* strictly necessary to have each component fully functional, but it can be quite useful.

For general information on serial console configuration under Linux, refer to the “Remote Serial Console HOWTO” at The Linux Documentation Project: <http://www.tldp.org>

Serial Console BIOS configuration

Enabling system serial console output neither enables nor disables serial capabilities in GRUB, Xen, or Linux, but may make remote management of your system more convenient by displaying POST and other boot messages over serial port and allowing remote BIOS configuration.

Refer to your hardware vendor’s documentation for capabilities and procedures to enable BIOS serial redirection.

Serial Console GRUB configuration

Enabling GRUB serial console output neither enables nor disables Xen or Linux serial capabilities, but may make remote management of your system more convenient by displaying GRUB prompts, menus, and actions over serial port and allowing remote GRUB management.

Adding the following two lines to your GRUB configuration file, typically either `/boot/grub/menu.lst` or `/boot/grub/grub.conf` depending on your distro, will enable GRUB serial output.

```
serial --unit=0 --speed=115200 --word=8 --parity=no --stop=1
terminal --timeout=10 serial console
```

Note that when both the serial port and the local monitor and keyboard are enabled, the text *“Press any key to continue”* will appear at both. Pressing a key on one device will cause GRUB to display to that device. The other device will see no output. If no key is pressed before the timeout period expires, the system will boot to the default GRUB boot entry.

Please refer to the GRUB documentation for further information.

Serial Console Xen configuration

Enabling Xen serial console output neither enables nor disables Linux kernel output or logging in to Linux over serial port. It does however allow you to monitor and log the Xen boot process via serial console and can be very useful in debugging.

In order to configure Xen serial console output, it is necessary to add a boot option to your GRUB config; e.g. replace the previous example kernel line with:

```
kernel /boot/xen.gz dom0_mem=131072 com1=115200,8n1
```

This configures Xen to output on COM1 at 115,200 baud, 8 data bits, no parity and 1 stop bit. Modify these parameters for your environment. See Section 11.3 for an explanation of all boot parameters.

One can also configure XenLinux to share the serial console; to achieve this append *“console=ttyS0”* to your module line.

Serial Console Linux configuration

Enabling Linux serial console output at boot neither enables nor disables logging in to Linux over serial port. It does however allow you to monitor and log the Linux boot process via serial console and can be very useful in debugging.

To enable Linux output at boot time, add the parameter `console=ttyS0` (or `ttyS1`, `ttyS2`, etc.) to your kernel GRUB line. Under Xen, this might be:

```
module /vmlinuz-2.6-xen0 ro root=/dev/VolGroup00/LogVol100 \
console=ttyS0, 115200
```

to enable output over `ttyS0` at 115200 baud.

Serial Console Login configuration

Logging in to Linux via serial console, under Xen or otherwise, requires specifying a login prompt be started on the serial port. To permit root logins over serial console, the serial port must be added to `/etc/securetty`.

To automatically start a login prompt over the serial port, add the line:

```
c:2345:respawn:/sbin/minigetty ttyS0
```

to `/etc/inittab`. Run `init q` to force a reload of your inittab and start getty.

To enable root logins, add `ttys0` to `/etc/securetty` if not already present.

Your distribution may use an alternate getty; options include `getty`, `mgetty` and `agetty`. Consult your distribution's documentation for further information.

2.5.3 TLS Libraries

Users of the XenLinux 2.6 kernel should disable Thread Local Storage (TLS) (e.g. by doing a `mv /lib/tls /lib/tls.disabled`) before attempting to boot a XenLinux kernel⁴. You can always reenable TLS by restoring the directory to its original location (i.e. `mv /lib/tls.disabled /lib/tls`).

The reason for this is that the current TLS implementation uses segmentation in a way that is not permissible under Xen. If TLS is not disabled, an emulation mode is used within Xen which reduces performance substantially. To ensure full performance you should install a 'Xen-friendly' (nosegneg) version of the library.

2.6 Booting Xen

It should now be possible to restart the system and use Xen. Reboot and choose the new Xen option when the Grub screen appears.

What follows should look much like a conventional Linux boot. The first portion of the output comes from Xen itself, supplying low level information about itself and the underlying hardware. The last portion of the output comes from XenLinux.

You may see some error messages during the XenLinux boot. These are not necessarily anything to worry about—they may result from kernel configuration differences between your XenLinux kernel and the one you usually use.

When the boot completes, you should be able to log into your system as usual. If you are unable to log in, you should still be able to reboot with your normal Linux kernel by selecting it at the GRUB prompt.

⁴If you boot without first disabling TLS, you will get a warning message during the boot process. In this case, simply perform the rename after the machine is up and then run `/sbin/ldconfig` to make it take effect.

Chapter 3

Booting a Xen System

Booting the system into Xen will bring you up into the privileged management domain, Domain0. At that point you are ready to create guest domains and “boot” them using the `xm create` command.

3.1 Booting Domain0

After installation and configuration is complete, reboot the system and choose the new Xen option when the Grub screen appears.

What follows should look much like a conventional Linux boot. The first portion of the output comes from Xen itself, supplying low level information about itself and the underlying hardware. The last portion of the output comes from XenLinux.

When the boot completes, you should be able to log into your system as usual. If you are unable to log in, you should still be able to reboot with your normal Linux kernel by selecting it at the GRUB prompt.

The first step in creating a new domain is to prepare a root filesystem for it to boot. Typically, this might be stored in a normal partition, an LVM or other volume manager partition, a disk file or on an NFS server. A simple way to do this is simply to boot from your standard OS install CD and install the distribution into another partition on your hard drive.

To start the xend control daemon, type

```
# xend start
```

If you wish the daemon to start automatically, see the instructions in Section 4.1. Once the daemon is running, you can use the `xm` tool to monitor and maintain the domains running on your system. This chapter provides only a brief tutorial. We provide full details of the `xm` tool in the next chapter.

3.2 Booting Guest Domains

3.2.1 Creating a Domain Configuration File

Before you can start an additional domain, you must create a configuration file. We provide two example files which you can use as a starting point:

- `/etc/xen/xmexample1` is a simple template configuration file for describing a single VM.
- `/etc/xen/xmexample2` file is a template description that is intended to be reused for multiple virtual machines. Setting the value of the `vmid` variable on the `xm` command line fills in parts of this template.

There are also a number of other examples which you may find useful. Copy one of these files and edit it as appropriate. Typical values you may wish to edit include:

kernel Set this to the path of the kernel you compiled for use with Xen
(e.g. `kernel = ``/boot/vmlinuz-2.6-xenU```)

memory Set this to the size of the domain's memory in megabytes (e.g.
`memory = 64`)

disk Set the first entry in this list to calculate the offset of the domain's root partition, based on the domain ID. Set the second to the location of `/usr` if you are sharing it between domains (e.g. `disk = ['phy:your_hard_drive%d,sda1,w' % (base_partition_number + vmid), 'phy:your_usr_partition,sda6,r']`)

dhcp Uncomment the `dhcp` variable, so that the domain will receive its IP address from a DHCP server (e.g. `dhcp=``dhcp```)

You may also want to edit the **vif** variable in order to choose the MAC address of the virtual ethernet interface yourself. For example:

```
vif = [ 'mac=00:16:3E:F6:BB:B3' ]
```

If you do not set this variable, `xend` will automatically generate a random MAC address from the range `00:16:3E:xx:xx:xx`, assigned by IEEE to XenSource as an OUI (organizationally unique identifier). XenSource Inc. gives permission for anyone to use addresses randomly allocated from this range for use by their Xen domains.

For a list of IEEE OUI assignments, see <http://standards.ieee.org/regauth/oui/oui.txt>

3.2.2 Booting the Guest Domain

The `xm` tool provides a variety of commands for managing domains. Use the `create` command to start new domains. Assuming you've created a configuration file `myvmconf`

based around `/etc/xen/xmexample2`, to start a domain with virtual machine ID 1 you should type:

```
# xm create -c myvmconf vmid=1
```

The `-c` switch causes `xm` to turn into the domain's console after creation. The `vmid=1` sets the `vmid` variable used in the `myvmconf` file.

You should see the console boot messages from the new domain appearing in the terminal in which you typed the command, culminating in a login prompt.

3.3 Starting / Stopping Domains Automatically

It is possible to have certain domains start automatically at boot time and to have `dom0` wait for all running domains to shutdown before it shuts down the system.

To specify a domain is to start at boot-time, place its configuration file (or a link to it) under `/etc/xen/auto/`.

A Sys-V style init script for Red Hat and LSB-compliant systems is provided and will be automatically copied to `/etc/init.d/` during install. You can then enable it in the appropriate way for your distribution.

For instance, on Red Hat:

```
# chkconfig --add xendomains
```

By default, this will start the boot-time domains in runlevels 3, 4 and 5.

You can also use the `service` command to run this script manually, e.g:

```
# service xendomains start
```

Starts all the domains with config files under `/etc/xen/auto/`.

```
# service xendomains stop
```

Shuts down all running Xen domains.

Part II

Configuration and Management

Chapter 4

Domain Management Tools

This chapter summarizes the management software and tools available.

4.1 Xend

The Xend node control daemon performs system management functions related to virtual machines. It forms a central point of control of virtualized resources, and must be running in order to start and manage virtual machines. Xend must be run as root because it needs access to privileged system management functions.

An initialization script named `/etc/init.d/xend` is provided to start Xend at boot time. Use the tool appropriate (i.e. `chkconfig`) for your Linux distribution to specify the runlevels at which this script should be executed, or manually create symbolic links in the correct runlevel directories.

Xend can be started on the command line as well, and supports the following set of parameters:

```
# xend start      start xend, if not already running
# xend stop       stop xend if already running
# xend restart    restart xend if running, otherwise start it
# xend status     indicates xend status by its return code
```

A SysV init script called `xend` is provided to start xend at boot time. `make install` installs this script in `/etc/init.d`. To enable it, you have to make symbolic links in the appropriate runlevel directories or use the `chkconfig` tool, where available. Once xend is running, administration can be done using the `xm` tool.

4.1.1 Logging

As xend runs, events will be logged to `/var/log/xen/xend.log` and (less frequently) to `/var/log/xen/xend-debug.log`. These, along with the standard syslog files, are useful when troubleshooting problems.

4.1.2 Configuring Xend

Xend is written in Python. At startup, it reads its configuration information from the file `/etc/xen/xend-config.sxp`. The Xen installation places an example `xend-config.sxp` file in the `/etc/xen` subdirectory which should work for most installations.

See the example configuration file `xend-debug.sxp` and the section 5 man page `xend-config.sxp` for a full list of parameters and more detailed information. Some of the most important parameters are discussed below.

An HTTP interface and a Unix domain socket API are available to communicate with Xend. This allows remote users to pass commands to the daemon. By default, Xend does not start an HTTP server. It does start a Unix domain socket management server, as the low level utility `xm` requires it. For support of cross-machine migration, Xend can start a relocation server. This support is not enabled by default for security reasons.

Note: the example xend configuration file modifies the defaults and starts up Xend as an HTTP server as well as a relocation server.

From the file:

```
 #(xend-http-server no)
 (xend-http-server yes)
 #(xend-unix-server yes)
 #(xend-relocation-server no)
 (xend-relocation-server yes)
```

Comment or uncomment lines in that file to disable or enable features that you require.

Connections from remote hosts are disabled by default:

```
 # Address xend should listen on for HTTP connections, if xend-http-server is
 # set.
 # Specifying 'localhost' prevents remote connections.
 # Specifying the empty string '' (the default) allows all connections.
 #(xend-address '')
 (xend-address localhost)
```

It is recommended that if migration support is not needed, the `xend-relocation-server` parameter value be changed to “no” or commented out.

4.2 Xm

The `xm` tool is the primary tool for managing Xen from the console. The general format of an `xm` command line is:

```
# xm command [switches] [arguments] [variables]
```

The available *switches* and *arguments* are dependent on the *command* chosen. The *variables* may be set using declarations of the form `variable=value` and command line declarations override any of the values in the configuration file being used, including the standard variables described above and any custom variables (for instance, the `xmdefconfig` file uses a `vmid` variable).

For online help for the commands available, type:

```
# xm help
```

This will list the most commonly used commands. The full list can be obtained using `xm help --long`. You can also type `xm help <command>` for more information on a given command.

4.2.1 Basic Management Commands

One useful command is `# xm list` which lists all domains running in rows of the following format:

```
name domid memory vcpus state cputime
```

The meaning of each field is as follows:

name The descriptive name of the virtual machine.

domid The number of the domain ID this virtual machine is running in.

memory Memory size in megabytes.

vcpus The number of virtual CPUs this domain has.

state Domain state consists of 5 fields:

r running

b blocked

p paused

s shutdown

c crashed

cputime How much CPU time (in seconds) the domain has used so far.

The `xm list` command also supports a long output format when the `-l` switch is used. This outputs the full details of the running domains in `xend`'s SXP configuration format.

If you want to know how long your domains have been running for, then you can use the `# xm uptime` command.

You can get access to the console of a particular domain using the `# xm console` command (e.g. `# xm console myVM`).

4.2.2 Domain Scheduling Management Commands

The credit CPU scheduler automatically load balances guest VCPUs across all available physical CPUs on an SMP host. The user need not manually pin VCPUs to load balance the system. However, she can restrict which CPUs a particular VCPU may run on using the `xm vcpu-pin` command.

Each guest domain is assigned a `weight` and a `cap`.

A domain with a weight of 512 will get twice as much CPU as a domain with a weight of 256 on a contended host. Legal weights range from 1 to 65535 and the default is 256.

The `cap` optionally fixes the maximum amount of CPU a guest will be able to consume, even if the host system has idle CPU cycles. The `cap` is expressed in percentage of one physical CPU: 100 is 1 physical CPU, 50 is half a CPU, 400 is 4 CPUs, etc... The default, 0, means there is no upper cap.

When you are running with the credit scheduler, you can check and modify your domains' weights and caps using the `xm sched-credit` command:

<code>xm sched-credit -d <domain></code>	lists weight and cap
<code>xm sched-credit -d <domain> -w <weight></code>	sets the weight
<code>xm sched-credit -d <domain> -c <cap></code>	sets the cap

Chapter 5

Domain Configuration

The following contains the syntax of the domain configuration files and description of how to further specify networking, driver domain and general scheduling behavior.

5.1 Configuration Files

Xen configuration files contain the following standard variables. Unless otherwise stated, configuration items should be enclosed in quotes: see the configuration scripts in `/etc/xen/` for concrete examples.

kernel Path to the kernel image.

ramdisk Path to a ramdisk image (optional).

memory Memory size in megabytes.

vcpus The number of virtual CPUs.

console Port to export the domain console on (default 9600 + domain ID).

vif Network interface configuration. This may simply contain an empty string for each desired interface, or may override various settings, e.g.

```
vif = [ 'mac=00:16:3E:00:00:11, bridge=xen-br0',  
        'bridge=xen-br1' ]
```

to assign a MAC address and bridge to the first interface and assign a different bridge to the second interface, leaving xen to choose the MAC address. The settings that may be overridden in this way are type, mac, bridge, ip, script, backend, and vifname.

disk List of block devices to export to the domain e.g. `disk = ['phy:hda1,sda1,r']` exports physical device `/dev/hda1` to the domain as `/dev/sda1` with read-only access. Exporting a disk read-write which is currently mounted is dangerous – if you are *certain* you wish to do this, you can specify `w!` as the mode.

dhcp Set to 'dhcp' if you want to use DHCP to configure networking.

netmask Manually configured IP netmask.

gateway Manually configured IP gateway.

hostname Set the hostname for the virtual machine.

root Specify the root device parameter on the kernel command line.

nfs_server IP address for the NFS server (if any).

nfs_root Path of the root filesystem on the NFS server (if any).

extra Extra string to append to the kernel command line (if any)

Additional fields are documented in the example configuration files (e.g. to configure virtual TPM functionality).

For additional flexibility, it is also possible to include Python scripting commands in configuration files. An example of this is the `xmexample2` file, which uses Python code to handle the `vmid` variable.

5.2 Network Configuration

For many users, the default installation should work “out of the box”. More complicated network setups, for instance with multiple Ethernet interfaces and/or existing bridging setups will require some special configuration.

The purpose of this section is to describe the mechanisms provided by xend to allow a flexible configuration for Xen’s virtual networking.

5.2.1 Xen virtual network topology

Each domain network interface is connected to a virtual network interface in dom0 by a point to point link (effectively a “virtual crossover cable”). These devices are named `vif<domid>.<vifid>` (e.g. `vif1.0` for the first interface in domain 1, `vif3.1` for the second interface in domain 3).

Traffic on these virtual interfaces is handled in domain 0 using standard Linux mechanisms for bridging, routing, rate limiting, etc. Xend calls on two shell scripts to perform initial configuration of the network and configuration of new virtual interfaces. By default, these scripts configure a single bridge for all the virtual interfaces. Arbitrary routing / bridging configurations can be configured by customizing the scripts, as described in the following section.

5.2.2 Xen networking scripts

Xen's virtual networking is configured by two shell scripts (by default `network-bridge` and `vif-bridge`). These are called automatically by `xend` when certain events occur, with arguments to the scripts providing further contextual information. These scripts are found by default in `/etc/xen/scripts`. The names and locations of the scripts can be configured in `/etc/xen/xend-config.sxp`.

network-bridge: This script is called whenever `xend` is started or stopped to respectively initialize or tear down the Xen virtual network. In the default configuration initialization creates the bridge 'xen-br0' and moves `eth0` onto that bridge, modifying the routing accordingly. When `xend` exits, it deletes the Xen bridge and removes `eth0`, restoring the normal IP and routing configuration.

vif-bridge: This script is called for every domain virtual interface and can configure firewalling rules and add the vif to the appropriate bridge. By default, this adds and removes VIFs on the default Xen bridge.

Other example scripts are available (`network-route` and `vif-route`, `network-nat` and `vif-nat`). For more complex network setups (e.g. where routing is required or integrate with existing bridges) these scripts may be replaced with customized variants for your site's preferred configuration.

5.3 Driver Domain Configuration

5.3.1 PCI

Individual PCI devices can be assigned to a given domain (a PCI driver domain) to allow that domain direct access to the PCI hardware.

While PCI Driver Domains can increase the stability and security of a system by addressing a number of security concerns, there are some security issues that remain that you can read about in Section 9.2.

Compile-Time Setup

To use this functionality, ensure that the PCI Backend is compiled in to a privileged domain (e.g. domain 0) and that the domains which will be assigned PCI devices have the PCI Frontend compiled in. In XenLinux, the PCI Backend is available under the Xen configuration section while the PCI Frontend is under the architecture-specific "Bus Options" section. You may compile both the backend and the frontend into the same kernel; they will not affect each other.

PCI Backend Configuration - Binding at Boot

The PCI devices you wish to assign to unprivileged domains must be "hidden" from your backend domain (usually domain 0) so that it does not load a driver for them. Use the `pciback.hide` kernel parameter which is specified on the kernel command-line and is configurable through GRUB (see Section 2.5). Note that devices are not really hidden from the backend domain. The PCI Backend appears to the Linux kernel as a regular PCI device driver. The PCI Backend ensures that no other device driver loads for the devices by binding itself as the device driver for those devices. PCI devices are identified by hexadecimal slot/function numbers (on Linux, use `lspci` to determine slot/function numbers of your devices) and can be specified with or without the PCI domain:

`(bus:slot.func) example (02:1d.3)`
`(domain:bus:slot.func) example (0000:02:1d.3)`

An example kernel command-line which hides two PCI devices might be:

```
root=/dev/sda4 ro console=tty0 pciback.hide=(02:01.f)(0000:04:1d.0)
```

PCI Backend Configuration - Late Binding

PCI devices can also be bound to the PCI Backend after boot through the manual binding/unbinding facilities provided by the Linux kernel in `sysfs` (allowing for a Xen user to give PCI devices to driver domains that were not specified on the kernel command-line). There are several attributes with the PCI Backend's `sysfs` directory (`/sys/bus/pci/drivers/pciback`) that can be used to bind/unbind devices:

slots lists all of the PCI slots that the PCI Backend will try to seize (or "hide" from Domain 0). A PCI slot must appear in this list before it can be bound to the PCI Backend through the `bind` attribute.

new_slot write the name of a slot here (in 0000:00:00.0 format) to have the PCI Backend seize the device in this slot.

remove_slot write the name of a slot here (same format as `new_slot`) to have the PCI Backend no longer try to seize devices in this slot. Note that this does not unbind the driver from a device it has already seized.

bind write the name of a slot here (in 0000:00:00.0 format) to have the Linux kernel attempt to bind the device in that slot to the PCI Backend driver.

unbind write the name of a slot here (same format as `bind`) to have the Linux kernel unbind the device from the PCI Backend. DO NOT unbind a device while it is currently given to a PCI driver domain!

Some examples:

Bind a device to the PCI Backend which is not bound to any other driver.

```
# # Add a new slot to the PCI Backend's list
# echo -n 0000:01:04.d > /sys/bus/pci/drivers/pciback/new_slot
# # Now that the backend is watching for the slot, bind to it
# echo -n 0000:01:04.d > /sys/bus/pci/drivers/pciback/bind
```

Unbind a device from its driver and bind to the PCI Backend.

```
# # Unbind a PCI network card from its network driver
# echo -n 0000:05:02.0 > /sys/bus/pci/drivers/3c905/unbind
# # And now bind it to the PCI Backend
# echo -n 0000:05:02.0 > /sys/bus/pci/drivers/pciback/new_slot
# echo -n 0000:05:02.0 > /sys/bus/pci/drivers/pciback/bind
```

Note that the "-n" option in the example is important as it causes echo to not output a new-line.

PCI Backend Configuration - User-space Quirks

Quirky devices (such as the Broadcom Tigon 3) may need write access to their configuration space registers. Xen can be instructed to allow specified PCI devices write access to specific configuration space registers. The policy may be found in:

```
/etc/xen/xend-pci-quirks.sxp
```

The policy file is heavily commented and is intended to provide enough documentation for developers to extend it.

PCI Backend Configuration - Permissive Flag

If the user-space quirks approach doesn't meet your needs you may want to enable the permissive flag for that device. To do so, first get the PCI domain, bus, slot, and function information from dom0 via `lspci`. Then augment the user-space policy for permissive devices. The permissive policy can be found in:

```
/etc/xen/xend-pci-permissive.sxp
```

Currently, the only way to reset the permissive flag is to unbind the device from the PCI Backend driver.

PCI Backend - Checking Status

There two important sysfs nodes that provide a mechanism to view specifics on quirks and permissive devices:

```
/sys/bus/drivers/pciback/permissive
```

Use `cat` on this file to view a list of permissive slots.

```
/sys/bus/drivers/pciback/quirks
```

Use `cat` on this file view a hierarchical view of devices bound to the PCI backend, their PCI vendor/device ID, and any quirks that are associated with that particular slot.

You may notice that every device bound to the PCI backend has 17 quirks standard "quirks" regardless of `xend-pci-quirks.sxp`. These default entries are necessary to support interactions between the PCI bus manager and the device bound to it. Even non-quirky devices should have these standard entries.

In this case, preference was given to accuracy over aesthetics by choosing to show the standard quirks in the quirks list rather than hide them from the inquiring user

PCI Frontend Configuration

To configure a domU to receive a PCI device:

Command-line: Use the `pci` command-line flag. For multiple devices, use the option multiple times.

```
xm create netcard-dd pci=01:00.0 pci=02:03.0
```

Flat Format configuration file: Specify all of your PCI devices in a python list named `pci`.

```
pci=[ '01:00.0' , '02:03.0' ]
```

SXP Format configuration file: Use a single PCI device section for all of your devices (specify the numbers in hexadecimal with the preceding '0x'). Note that *domain* here refers to the PCI domain, not a virtual machine within Xen.

```
(device (pci
  (dev (domain 0x0)(bus 0x3)(slot 0x1a)(func 0x1)
  (dev (domain 0x0)(bus 0x1)(slot 0x5)(func 0x0)
)
```

5.4 Support for virtual Trusted Platform Module (vTPM)

Paravirtualized domains can be given access to a virtualized version of a TPM. This enables applications in these domains to use the services of the TPM device for example through a TSS stack ¹. The Xen source repository provides the necessary software components to enable virtual TPM access. Support is provided through several different pieces. First, a TPM emulator has been modified to provide TPM's functionality for the virtual TPM subsystem. Second, a virtual TPM Manager coordinates the virtual TPMs efforts, manages their creation, and provides protected key storage using

¹Trousers TSS stack: <http://sourceforge.net/projects/trousers>

the TPM. Third, a device driver pair providing a TPM front- and backend is available for XenLinux to deliver TPM commands from the domain to the virtual TPM manager, which dispatches it to a software TPM. Since the TPM Manager relies on a HW TPM for protected key storage, therefore this subsystem requires a Linux-supported hardware TPM. For development purposes, a TPM emulator is available for use on non-TPM enabled platforms.

Compile-Time Setup

To enable access to the virtual TPM, the virtual TPM backend driver must be compiled for a privileged domain (e.g. domain 0). Using the XenLinux configuration, the necessary driver can be selected in the Xen configuration section. Unless the driver has been compiled into the kernel, its module must be activated using the following command:

```
modprobe tpmbk
```

Similarly, the TPM frontend driver must be compiled for the kernel trying to use TPM functionality. Its driver can be selected in the kernel configuration section Device Driver / Character Devices / TPM Devices. Along with that the TPM driver for the built-in TPM must be selected. If the virtual TPM driver has been compiled as module, it must be activated using the following command:

```
modprobe tpm_xenu
```

Furthermore, it is necessary to build the virtual TPM manager and software TPM by making changes to entries in Xen build configuration files. The following entry in the file Config.mk in the Xen root source directory must be made:

```
VTPM_TOOLS ?= y
```

After a build of the Xen tree and a reboot of the machine, the TPM backend drive must be loaded. Once loaded, the virtual TPM manager daemon must be started before TPM-enabled guest domains may be launched. To enable being the destination of a virtual TPM Migration, the virtual TPM migration daemon must also be loaded.

```
vtpm_managerd
```

```
vtpm_migratord
```

Once the VTPM manager is running, the VTPM can be accessed by loading the front end driver in a guest domain.

Development and Testing TPM Emulator

For development and testing on non-TPM enabled platforms, a TPM emulator can be used in replacement of a platform TPM. First, the entry in the file tools/vtpm/Rules.mk must look as follows:

```
BUILD_EMULATOR = y
```

Second, the entry in the file `tool/vtpm_manager/Rules.mk` must be uncommented as follows:

```
# TCS talks to fifo's rather than /dev/tpm. TPM Emulator assumed on fifos
CFLAGS += -DDUMMY_TPM
```

Before starting the virtual TPM Manager, start the emulator by executing the following in `dom0`:

```
tpm_emulator clear
```

vTPM Frontend Configuration

To provide TPM functionality to a user domain, a line must be added to the virtual TPM configuration file using the following format:

```
vtpm = ['instance=<instance number>, backend=<domain id>']
```

The *instance number* reflects the preferred virtual TPM instance to associate with the domain. If the selected instance is already associated with another domain, the system will automatically select the next available instance. An instance number greater than zero must be provided. It is possible to omit the instance parameter from the configuration file.

The *domain id* provides the ID of the domain where the virtual TPM backend driver and virtual TPM are running in. It should currently always be set to '0'.

Examples for valid vtpm entries in the configuration file are

```
vtpm = ['instance=1, backend=0']
```

and

```
vtpm = ['backend=0'].
```

Using the virtual TPM

Access to TPM functionality is provided by the virtual TPM frontend driver. Similar to existing hardware TPM drivers, this driver provides basic TPM status information through the *sysfs* filesystem. In a Xen user domain the *sysfs* entries can be found in `/sys/devices/xen/vtpm-0`.

Commands can be sent to the virtual TPM instance using the character device `/dev/tpm0` (major 10, minor 224).

Chapter 6

Storage and File System Management

Storage can be made available to virtual machines in a number of different ways. This chapter covers some possible configurations.

The most straightforward method is to export a physical block device (a hard drive or partition) from dom0 directly to the guest domain as a virtual block device (VBD).

Storage may also be exported from a filesystem image or a partitioned filesystem image as a *file-backed VBD*.

Finally, standard network storage protocols such as NBD, iSCSI, NFS, etc., can be used to provide storage to virtual machines.

6.1 Exporting Physical Devices as VBDs

One of the simplest configurations is to directly export individual partitions from domain 0 to other domains. To achieve this use the `phy:` specifier in your domain configuration file. For example a line like

```
disk = [ 'phy:hda3,sda1,w' ]
```

specifies that the partition `/dev/hda3` in domain 0 should be exported read-write to the new domain as `/dev/sda1`; one could equally well export it as `/dev/hda` or `/dev/sdb5` should one wish.

In addition to local disks and partitions, it is possible to export any device that Linux considers to be “a disk” in the same manner. For example, if you have iSCSI disks or GNBD volumes imported into domain 0 you can export these to other domains using the `phy:` disk syntax. E.g.:

```
disk = [ 'phy:vg/lvm1,sda2,w' ]
```

Warning: Block device sharing

Block devices should typically only be shared between domains in a read-only fashion otherwise the Linux kernel's file systems will get very confused as the file system structure may change underneath them (having the same ext3 partition mounted `rw` twice is a sure fire way to cause irreparable damage)! Xend will attempt to prevent you from doing this by checking that the device is not mounted read-write in domain 0, and hasn't already been exported read-write to another domain. If you want read-write sharing, export the directory to other domains via NFS from domain 0 (or use a cluster file system such as GFS or ocfs2).

6.2 Using File-backed VBDs

It is also possible to use a file in Domain 0 as the primary storage for a virtual machine. As well as being convenient, this also has the advantage that the virtual block device will be *sparse* — space will only really be allocated as parts of the file are used. So if a virtual machine uses only half of its disk space then the file really takes up half of the size allocated.

For example, to create a 2GB sparse file-backed virtual block device (actually only consumes 1KB of disk):

```
# dd if=/dev/zero of=vmldisk bs=1k seek=2048k count=1
```

Make a file system in the disk file:

```
# mkfs -t ext3 vmldisk
```

(when the tool asks for confirmation, answer 'y')

Populate the file system e.g. by copying from the current root:

```
# mount -o loop vmldisk /mnt
# cp -ax /{root,dev,var,etc,usr,bin,sbin,lib} /mnt
# mkdir /mnt/{proc,sys,home,tmp}
```

Tailor the file system by editing `/etc/fstab`, `/etc/hostname`, etc. Don't forget to edit the files in the mounted file system, instead of your domain 0 filesystem, e.g. you would edit `/mnt/etc/fstab` instead of `/etc/fstab`. For this example put `/dev/sda1` to root in `fstab`.

Now unmount (this is important!):

```
# umount /mnt
```

In the configuration file set:

```
disk = [ 'tap:aio:/full/path/to/vmldisk,sda1,w' ]
```


As the virtual machine writes to its ‘disk’, the sparse file will be filled in and consume more space up to the original 2GB.

Note: Users that have worked with file-backed VBDs on Xen in previous versions will be interested to know that this support is now provided through the blktap driver instead of the loopback driver. This change results in file-based block devices that are higher-performance, more scalable, and which provide better safety properties for VBD data. All that is required to update your existing file-backed VM configurations is to change VBD configuration lines from:

```
disk = [ 'file:/full/path/to/vmldisk,sda1,w' ]
```

to:

```
disk = [ 'tap:aio:/full/path/to/vmldisk,sda1,w' ]
```

6.2.1 Loopback-mounted file-backed VBDs (deprecated)

Note: *Loopback mounted VBDs have now been replaced with blktap-based support for raw image files, as described above. This section remains to detail a configuration that was used by older Xen versions.*

Raw image file-backed VBDs may also be attached to VMs using the Linux loopback driver. The only required change to the raw file instructions above are to specify the configuration entry as:

```
disk = [ 'file:/full/path/to/vmldisk,sda1,w' ]
```

Note that loopback file-backed VBDs may not be appropriate for backing I/O-intensive domains. This approach is known to experience substantial slowdowns under heavy I/O workloads, due to the I/O handling by the loopback block device used to support file-backed VBDs in dom0. Loopback support remains for old Xen installations, and users are strongly encouraged to use the blktap-based file support (using “tap:aio” as described above).

Additionally, Linux supports a maximum of eight loopback file-backed VBDs across all domains by default. This limit can be statically increased by using the *max_loop* module parameter if CONFIG_BLK_DEV_LOOP is compiled as a module in the dom0 kernel, or by using the *max_loop=n* boot option if CONFIG_BLK_DEV_LOOP is compiled directly into the dom0 kernel. Again, users are encouraged to use the blktap-based file support described above which scales to much larger number of active VBDs.

6.3 Using LVM-backed VBDs

A particularly appealing solution is to use LVM volumes as backing for domain file-systems since this allows dynamic growing/shrinking of volumes as well as snapshot and other features.

To initialize a partition to support LVM volumes:

```
# pvcreate /dev/sda10
```

Create a volume group named 'vg' on the physical partition:

```
# vgcreate vg /dev/sda10
```

Create a logical volume of size 4GB named 'myvmdisk1':

```
# lvcreate -L4096M -n myvmdisk1 vg
```

You should now see that you have a /dev/vg/myvmdisk1. Make a filesystem, mount it and populate it, e.g.:

```
# mkfs -t ext3 /dev/vg/myvmdisk1
# mount /dev/vg/myvmdisk1 /mnt
# cp -ax / /mnt
# umount /mnt
```

Now configure your VM with the following disk configuration:

```
disk = [ 'phy:vg/myvmdisk1,sda1,w' ]
```

LVM enables you to grow the size of logical volumes, but you'll need to resize the corresponding file system to make use of the new space. Some file systems (e.g. ext3) now support online resize. See the LVM manuals for more details.

You can also use LVM for creating copy-on-write (CoW) clones of LVM volumes (known as writable persistent snapshots in LVM terminology). This facility is new in Linux 2.6.8, so isn't as stable as one might hope. In particular, using lots of CoW LVM disks consumes a lot of dom0 memory, and error conditions such as running out of disk space are not handled well. Hopefully this will improve in future.

To create two copy-on-write clones of the above file system you would use the following commands:

```
# lvcreate -s -L1024M -n myclonedisk1 /dev/vg/myvmdisk1
# lvcreate -s -L1024M -n myclonedisk2 /dev/vg/myvmdisk1
```

Each of these can grow to have 1GB of differences from the master volume. You can grow the amount of space for storing the differences using the `lvextend` command, e.g.:

```
# lvextend +100M /dev/vg/myclonedisk1
```

Don't let the 'differences volume' ever fill up otherwise LVM gets rather confused. It may be possible to automate the growing process by using `dmsetup wait` to spot the volume getting full and then issue an `lvextend`.

In principle, it is possible to continue writing to the volume that has been cloned (the changes will not be visible to the clones), but we wouldn't recommend this: have the cloned volume as a 'pristine' file system install that isn't mounted directly by any of the virtual machines.

6.4 Using NFS Root

First, populate a root filesystem in a directory on the server machine. This can be on a distinct physical machine, or simply run within a virtual machine on the same node.

Now configure the NFS server to export this filesystem over the network by adding a line to `/etc/exports`, for instance:

```
/export/vmlroot      1.2.3.4/24 (rw,sync,no_root_squash)
```

Finally, configure the domain to use NFS root. In addition to the normal variables, you should make sure to set the following values in the domain's configuration file:

```
root      = '/dev/nfs'
nfs_server = '2.3.4.5'      # substitute IP address of server
nfs_root   = '/path/to/root' # path to root FS on the server
```

The domain will need network access at boot time, so either statically configure an IP address using the config variables `ip`, `netmask`, `gateway`, `hostname`; or enable DHCP (`dhcp='dhcp'`).

Note that the Linux NFS root implementation is known to have stability problems under high load (this is not a Xen-specific problem), so this configuration may not be appropriate for critical servers.

Chapter 7

CPU Management

Xen allows a domain's virtual CPU(s) to be associated with one or more host CPUs. This can be used to allocate real resources among one or more guests, or to make optimal use of processor resources when utilizing dual-core, hyperthreading, or other advanced CPU technologies.

Xen enumerates physical CPUs in a 'depth first' fashion. For a system with both hyperthreading and multiple cores, this would be all the hyperthreads on a given core, then all the cores on a given socket, and then all sockets. I.e. if you had a two socket, dual core, hyperthreaded Xeon the CPU order would be:

socket0				socket1			
core0		core1		core0		core1	
ht0	ht1	ht0	ht1	ht0	ht1	ht0	ht1
#0	#1	#2	#3	#4	#5	#6	#7

Having multiple vcpus belonging to the same domain mapped to the same physical CPU is very likely to lead to poor performance. It's better to use 'vcpu-set' to hot-unplug one of the vcpus and ensure the others are pinned on different CPUs.

If you are running IO intensive tasks, its typically better to dedicate either a hyper-thread or whole core to running domain 0, and hence pin other domains so that they can't use CPU 0. If your workload is mostly compute intensive, you may want to pin vcpus such that all physical CPU threads are available for guest domains.

Chapter 8

Migrating Domains

8.1 Domain Save and Restore

The administrator of a Xen system may suspend a virtual machine's current state into a disk file in domain 0, allowing it to be resumed at a later time.

For example you can suspend a domain called “VM1” to disk using the command:

```
# xm save VM1 VM1.chk
```

This will stop the domain named “VM1” and save its current state into a file called `VM1.chk`.

To resume execution of this domain, use the `xm restore` command:

```
# xm restore VM1.chk
```

This will restore the state of the domain and resume its execution. The domain will carry on as before and the console may be reconnected using the `xm console` command, as described earlier.

8.2 Migration and Live Migration

Migration is used to transfer a domain between physical hosts. There are two varieties: regular and live migration. The former moves a virtual machine from one host to another by pausing it, copying its memory contents, and then resuming it on the destination. The latter performs the same logical functionality but without needing to pause the domain for the duration. In general when performing live migration the domain continues its usual activities and—from the user's perspective—the migration should be imperceptible.

To perform a live migration, both hosts must be running Xen / `xend` and the destination host must have sufficient resources (e.g. memory capacity) to accommodate the

domain after the move. Furthermore we currently require both source and destination machines to be on the same L2 subnet.

Currently, there is no support for providing automatic remote access to filesystems stored on local disk when a domain is migrated. Administrators should choose an appropriate storage solution (i.e. SAN, NAS, etc.) to ensure that domain filesystems are also available on their destination node. GNBD is a good method for exporting a volume from one machine to another. iSCSI can do a similar job, but is more complex to set up.

When a domain migrates, it's MAC and IP address move with it, thus it is only possible to migrate VMs within the same layer-2 network and IP subnet. If the destination node is on a different subnet, the administrator would need to manually configure a suitable etherip or IP tunnel in the domain 0 of the remote node.

A domain may be migrated using the `xm migrate` command. To live migrate a domain to another machine, we would use the command:

```
# xm migrate --live mydomain destination.ournetwork.com
```

Without the `--live` flag, `xend` simply stops the domain and copies the memory image over to the new node and restarts it. Since domains can have large allocations this can be quite time consuming, even on a Gigabit network. With the `--live` flag `xend` attempts to keep the domain running while the migration is in progress, resulting in typical down times of just 60–300ms.

For now it will be necessary to reconnect to the domain's console on the new machine using the `xm console` command. If a migrated domain has any open network connections then they will be preserved, so SSH connections do not have this limitation.

Chapter 9

Securing Xen

This chapter describes how to secure a Xen system. It describes a number of scenarios and provides a corresponding set of best practices. It begins with a section devoted to understanding the security implications of a Xen system.

9.1 Xen Security Considerations

When deploying a Xen system, one must be sure to secure the management domain (Domain-0) as much as possible. If the management domain is compromised, all other domains are also vulnerable. The following are a set of best practices for Domain-0:

1. **Run the smallest number of necessary services.** The less things that are present in a management partition, the better. Remember, a service running as root in the management domain has full access to all other domains on the system.
2. **Use a firewall to restrict the traffic to the management domain.** A firewall with default-reject rules will help prevent attacks on the management domain.
3. **Do not allow users to access Domain-0.** The Linux kernel has been known to have local-user root exploits. If you allow normal users to access Domain-0 (even as unprivileged users) you run the risk of a kernel exploit making all of your domains vulnerable.

9.2 Driver Domain Security Considerations

Driver domains address a range of security problems that exist regarding the use of device drivers and hardware. On many operating systems in common use today, device drivers run within the kernel with the same privileges as the kernel. Few or no mechanisms exist to protect the integrity of the kernel from a misbehaving (read "buggy") or

malicious device driver. Driver domains exist to aid in isolating a device driver within its own virtual machine where it cannot affect the stability and integrity of other domains. If a driver crashes, the driver domain can be restarted rather than have the entire machine crash (and restart) with it. Drivers written by unknown or untrusted third-parties can be confined to an isolated space. Driver domains thus address a number of security and stability issues with device drivers.

However, due to limitations in current hardware, a number of security concerns remain that need to be considered when setting up driver domains (it should be noted that the following list is not intended to be exhaustive).

1. **Without an IOMMU, a hardware device can DMA to memory regions outside of its controlling domain.** Architectures which do not have an IOMMU (e.g. most x86-based platforms) to restrict DMA usage by hardware are vulnerable. A hardware device which can perform arbitrary memory reads and writes can read/write outside of the memory of its controlling domain. A malicious or misbehaving domain could use a hardware device it controls to send data overwriting memory in another domain or to read arbitrary regions of memory in another domain.
2. **Shared buses are vulnerable to sniffing.** Devices that share a data bus can sniff (and possibly spoof) each others' data. Device A that is assigned to Domain A could eavesdrop on data being transmitted by Domain B to Device B and then relay that data back to Domain A.
3. **Devices which share interrupt lines can either prevent the reception of that interrupt by the driver domain or can trigger the interrupt service routine of that guest needlessly.** A devices which shares a level-triggered interrupt (e.g. PCI devices) with another device can raise an interrupt and never clear it. This effectively blocks other devices which share that interrupt line from notifying their controlling driver domains that they need to be serviced. A device which shares an any type of interrupt line can trigger its interrupt continually which forces execution time to be spent (in multiple guests) in the interrupt service routine (potentially denying time to other processes within that guest). System architectures which allow each device to have its own interrupt line (e.g. PCI's Message Signaled Interrupts) are less vulnerable to this denial-of-service problem.
4. **Devices may share the use of I/O memory address space.** Xen can only restrict access to a device's physical I/O resources at a certain granularity. For interrupt lines and I/O port address space, that granularity is very fine (per interrupt line and per I/O port). However, Xen can only restrict access to I/O memory address space on a page size basis. If more than one device shares use of a page in I/O memory address space, the domains to which those devices are assigned will be able to access the I/O memory address space of each other's devices.

9.3 Security Scenarios

9.3.1 The Isolated Management Network

In this scenario, each node has two network cards in the cluster. One network card is connected to the outside world and one network card is a physically isolated management network specifically for Xen instances to use.

As long as all of the management partitions are trusted equally, this is the most secure scenario. No additional configuration is needed other than forcing Xend to bind to the management interface for relocation.

9.3.2 A Subnet Behind a Firewall

In this scenario, each node has only one network card but the entire cluster sits behind a firewall. This firewall should do at least the following:

1. Prevent IP spoofing from outside of the subnet.
2. Prevent access to the relocation port of any of the nodes in the cluster except from within the cluster.

The following iptables rules can be used on each node to prevent migrations to that node from outside the subnet assuming the main firewall does not do this for you:

```
# this command disables all access to the Xen relocation
# port:
iptables -A INPUT -p tcp --destination-port 8002 -j REJECT

# this command enables Xen relocations only from the specific
# subnet:
iptables -I INPUT -p tcp -{}-source 192.168.1.1/8 \
    --destination-port 8002 -j ACCEPT
```

9.3.3 Nodes on an Untrusted Subnet

Migration on an untrusted subnet is not safe in current versions of Xen. It may be possible to perform migrations through a secure tunnel via an VPN or SSH. The only safe option in the absence of a secure tunnel is to disable migration completely. The easiest way to do this is with iptables:

```
# this command disables all access to the Xen relocation port
iptables -A INPUT -p tcp -{}-destination-port 8002 -j REJECT
```


Chapter 10

sHype/Xen Access Control

The Xen mandatory access control framework is an implementation of the sHype Hypervisor Security Architecture (www.research.ibm.com/ssd_shype). It permits or denies communication and resource access of domains based on a security policy. The mandatory access controls are enforced in addition to the Xen core controls, such as memory protection. They are designed to remain transparent during normal operation of domains (policy-conform behavior) but to intervene when domains move outside their intended sharing behavior. This chapter will describe how the sHype access controls in Xen can be configured to prevent viruses from spilling over from one into another workload type and secrets from leaking from one workload type to another. sHype/Xen depends on the correct behavior of Domain0 (cf previous chapter).

Benefits of configuring sHype/ACM in Xen include:

- robust workload and resource protection effective against rogue user domains
- simple, platform- and operating system-independent security policies (ideal for heterogeneous distributed environments)
- safety net with minimal performance overhead in case operating system security is missing, does not scale, or fails

These benefits are very valuable because today's operating systems become increasingly complex and often have no or insufficient mandatory access controls. (Discretionary access controls, supported by most operating systems, are not effective against viruses or misbehaving programs.) Where mandatory access control exists (e.g., SELinux), they usually deploy complex and difficult to understand security policies. Additionally, multi-tier applications in business environments usually require different types of operating systems (e.g., AIX, Windows, Linux) which cannot be configured with compatible security policies. Related distributed transactions and workloads cannot be easily protected on the OS level. The Xen access control framework steps in to offer a coarse-grained but very robust security layer and safety net in case operating system security fails or is missing.



Figure 10.1: Overview of activating sHype workload protection in Xen. Section numbers point to representative examples.

To control sharing between domains, Xen mediates all inter-domain communication (shared memory, events) as well as the access of domains to resources such as disks. Thus, Xen can confine distributed workloads (domain payloads) by permitting sharing among domains running the same type of workload and denying sharing between pairs of domains that run different workload types. We assume that—from a Xen perspective—only one workload type is running per user domain. To enable Xen to associate domains and resources with workload types, security labels including the workload types are attached to domains and resources. These labels and the hypervisor sHype controls cannot be manipulated or bypassed and are effective even against rogue domains.

10.1 Overview

This section gives an overview of how workloads can be protected using the sHype mandatory access control framework in Xen. Figure 10.1 shows the necessary steps in activating the Xen workload protection. These steps are described in detail in Section 10.2.

First, the sHype/ACM access control must be enabled in the Xen distribution and the distribution must be built and installed (cf Subsection 10.2.1). Before we can enforce security, a Xen security policy must be created (cf Subsection 10.2.2) and deployed (cf Subsection 10.2.3). This policy defines the workload types differentiated during access control. It also defines the rules that compare workload types of domains and resources to provide access decisions. Workload types are represented by security labels that can be attached to domains and resources (cf Subsections 10.2.4 and 10.2.5). The

functioning of the active sHype/Xen workload protection is demonstrated using simple resource assignment, and domain creation tests in Subsection 10.2.6. Section 10.3 describes the syntax and semantics of the sHype/Xen security policy in detail and introduces briefly the tools that are available to help create valid security policies.

The next section describes all the necessary steps to create, deploy, and test a simple workload protection policy. It is meant to enable anybody to quickly try out the sHype/Xen workload protection. Those readers who are interested in learning more about how the sHype access control in Xen works and how it is configured using the XML security policy should read Section 10.3 as well. Section 10.4 concludes this chapter with current limitations of the sHype implementation for Xen.

10.2 Xen Workload Protection Step-by-Step

What you are about to do consists of the following sequence:

- configure and install sHype/Xen
- create a simple workload protection security policy
- deploy the sHype/Xen security policy
- associate domains and resources with workload labels,
- test the workload protection

The essential commands to create and deploy a sHype/Xen security policy are numbered throughout the following sections. If you want a quick-guide or return at a later time to go quickly through this demonstration, simply look for the numbered commands and apply them in order.

10.2.1 Configuring/Building sHype Support into Xen

First, we need to configure the access control module in Xen and install the ACM-enabled Xen hypervisor. This step installs security tools and compiles sHype/ACM controls into the Xen hypervisor.

To enable sHype/ACM in Xen, please edit the Config.mk file in the top Xen directory.

```
(1) In Config.mk
    Change: ACM_SECURITY ?= n
    To:    ACM_SECURITY ?= y
```

Then install the security-enabled Xen environment as follows:

```
(2) # make world
    # make install
```

10.2.2 Creating A WLP Policy in 3 Simple Steps with ezPolicy

We will use the ezPolicy tool to quickly create a policy that protects workloads. You will need both the Python and wxPython packages to run this tool. To run the tool in Domain0, you can download the wxPython package from www.wxpython.org or use the command `yum install wxPython` in Redhat/Fedora. To run the tool on MS Windows, you also need to download the Python package from www.python.org. After these packages are installed, start the ezPolicy tool with the following command:

```
(3) # xensec_ezpolicy
```

Figure 10.2 shows a screen-shot of the tool. The following steps show you how to create the policy shown in Figure 10.2. You can use <CTRL>-h to pop up a help window at any time. The indicators (a), (b), and (c) in Figure 10.2 show the buttons that are used during the 3 steps of creating a policy:

1. defining workloads
2. defining run-time conflicts
3. translating the workload definition into a sHype/Xen access control policy

Defining workloads. Workloads are defined for each organization and department that you enter in the left panel. Please use the “New Org” button (a) to create the organizations “Avis”, “Hertz”, “CocaCola”, and “PepsiCo”.

You can refine an organization to differentiate between multiple department workloads by right-clicking the organization and selecting Add Department (or selecting an organization and pressing <CTRL>-a). Create department workloads “Intranet”, “Extranet”, “HumanResources”, and “Payroll” for the “CocaCola” organization and department workloads “Intranet” and “Extranet” for the “PepsiCo” organization. The resulting layout of the tool should be similar to the left panel shown in Figure 10.2.

Defining run-time conflicts. Workloads that shall be prohibited from running concurrently on the same hypervisor platform are grouped into “Run-time Exclusion rules” on the right panel of the window.

To prevent PepsiCo and CocaCola workloads (including their departmental workloads) from running simultaneously on the same hypervisor system, select the organization “PepsiCo” and, while pressing the <CTRL>-key, select the organization “CocaCola”. Now press the button (b) named “Create run-time exclusion rule from selection”. A popup window will ask for the name for this run-time exclusion rule (enter a name or just hit <ENTER>). A rule will appear on the right panel. The name is used as reference only and does not affect the hypervisor policy.

Repeat the process to create a run-time exclusion rule just for the department workloads CocaCola.Extranet and CocaCola.Payroll.



Figure 10.2: Final layout including workload definition and Run-time Exclusion rules.

The resulting layout of your window should be similar to Figure 10.2. Save this workload definition by selecting “Save Workload Definition as ...” in the “File” menu (c). This workload definition can be later refined if required.

Translating the workload definition into a sHype/Xen access control policy. To translate the workload definition into a access control policy understood by Xen, please select the “Save as Xen ACM Security Policy” in the “File” menu (c). Enter the following policy name in the popup window: `example.chwall_ste.test-wld`. If you are running ezPolicy in Domain0, the resulting policy file `test-wld.security-policy.xml` will automatically be placed into the right directory (`/etc/xen/acm-security/policies/example/chwall_ste`). If you run the tool on another system, then you need to copy the resulting policy file into Domain0 before continuing. See Section 10.3.1 for naming conventions of security policies.

10.2.3 Deploying a WLP Policy

To deploy the workload protection policy we created in Section 10.2.2, we create a policy representation (`test-wld.bin`) that can be loaded into the Xen hypervisor and we configure Xen to actually load this policy at startup time.

The following command translates the source policy representation into a format that can be loaded into Xen with sHype/ACM support. Refer to the `xm` man page for further details:

```
(4) # xm makepolicy example.chwall_ste.test-wld
```

The easiest way to install a security policy for Xen is to include the policy in the boot sequence. The following command does just this:

```
(5) # xm cfgbootpolicy example.chwall_ste.test-wld
```

Alternatively, if this command fails (e.g., because it cannot identify the Xen boot entry), you can manually install the policy in 2 steps. First, manually copy the policy binary file into the boot directory:

```
# cp /etc/xen/acm-security/policies/example/chwall_ste/test-wld.bin \
/boot/example.chwall_ste.test-wld.bin
```

Second, manually add a module line to your Xen boot entry so that grub loads this policy file during startup:

```
title Xen (2.6.16.13)
    root (hd0,0)
    kernel /xen.gz dom0_mem=2000000 console=vga
    module /vmlinuz-2.6.16.13-xen ro root=/dev/hda3
    module /initrd-2.6.16.13-xen.img
    module /example.chwall_ste.test-wld.bin
```

Now reboot into this Xen boot entry to activate the policy and the security-enabled Xen hypervisor.

```
(6) # reboot
```

After reboot, check if security is enabled:

```
# xm list --label
Name          ID Mem(MiB) VCPUs State  Time(s)  Label
Domain-0      0    1949      4 r----- 163.9    SystemManagement
```

If the security label at the end of the line says “INACTIV” then the security is not enabled. Verify the previous steps. Note: Domain0 is assigned a default label (see `bootstrap` policy attribute explained in Section 10.3). All other domains must be labeled in order to start on this sHype/ACM-enabled Xen hypervisor (see following sections for labeling domains and resources).

10.2.4 Labeling Domains

You should have a Xen domain configuration file that looks like the following (Note: www.jailtime.org or www.xen-get.org might be good places to look for example domU images). The following configuration file defines `domain1`:

```
# cat domain1.xm
kernel = "/boot/vmlinuz-2.6.16.13-xen"
memory = 128
name = "domain1"
vif = [ '' ]
dhcp = "dhcp"
disk = [ 'file:/home/xen/dom_fc5/fedora.fc5.img,sda1,w', \
        'file:/home/xen/dom_fc5/fedora.swap,sda2,w' ]
root = "/dev/sda1 ro"
```

If you try to start domain1, you will get the following error:

```
# xm create domain1.xml
Using config file "domain1.xml".
domain1: DENIED
--> Domain not labeled
Checking resources: (skipped)
Security configuration prevents domain from starting
```

Every domain must be associated with a security label before it can start on sHype/Xen. Otherwise, sHype/Xen would not be able to enforce the policy consistently. The following command prints all domain labels available in the active policy:

```
# xm labels type=dom
Avis
CocaCola
CocaCola.Extranet
CocaCola.HumanResources
CocaCola.Intranet
CocaCola.Payroll
Hertz
PepsiCo
PepsiCo.Extranet
PepsiCo.Intranet
SystemManagement
```

Now label domain1 with the CocaCola label and another domain2 with the PepsiCo.Extranet label. Please refer to the xm man page for further information.

```
(7) # xm addlabel CocaCola dom domain1.xml
# xm addlabel PepsiCo.Extranet dom domain2.xml
```

Let us try to start the domain again:

```
# xm create domain1.xml
Using config file "domain1.xml".
file:/home/xen/dom_fc5/fedora.fc5.img: DENIED
--> res:__NULL_LABEL__ (NULL)
--> dom:CocaCola (example.chwall_ste.test-wld)
file:/home/xen/dom_fc5/fedora.swap: DENIED
--> res:__NULL_LABEL__ (NULL)
--> dom:CocaCola (example.chwall_ste.test-wld)
Security configuration prevents domain from starting
```

This error indicates that domain1, if started, would not be able to access its image and swap files because they are not labeled. This makes sense because to confine workloads, access of domains to resources must be controlled. Otherwise, domains that are not allowed to communicate or run simultaneously could share data through storage resources.

10.2.5 Labeling Resources

You can use the `xm labels type=res` command to list available resource labels. Let us assign the CocaCola resource label to the domain1 image file representing `/dev/sda1` and to its swap file:

```
(8) # xm addlabel CocaCola res \
      file:/home/xen/dom_fc5/fedora.fc5.img
Resource file not found, creating new file at:
/etc/xen/acm-security/policies/resource_labels
# xm addlabel CocaCola res \
      file:/home/xen/dom_fc5/fedora.swap
```

Starting domain1 now will succeed:

```
# xm create domain1.xm
# xm list --label
```

Name	ID	Mem(MiB)	VCPUs	State	Time(s)	Label
domain1	1	128	1	r-----	2.8	CocaCola
Domain-0	0	1949	4	r-----	387.7	SystemManagement

The following command lists all labeled resources on the system, e.g., to lookup or verify the labeling:

```
# xm resources
file:/home/xen/dom_fc5/fedora.swap
policy: example.chwall_ste.test-wld
label: CocaCola
file:/home/xen/dom_fc5/fedora.fc5.img
policy: example.chwall_ste.test-wld
label: CocaCola
```

Currently, if a labeled resource is moved to another location, the label must first be manually removed, and after the move re-attached using the `xm rmlabel` and `xm addlabel` respectively. Please see Section 10.4 for further details.

```
(9) Label the resources of domain2 as PepsiCo.Extranet
Do not try to start this domain yet
```

10.2.6 Testing The Xen Workload Protection

We are about to demonstrate how the workload protection works by verifying:

- that domains with conflicting workloads cannot run simultaneously
- that domains cannot access resources of other workloads
- that domains cannot exchange network packets if they are not associated with the same workload type

Test 1: Run-time exclusion rules. We assume that domain1 with the CocaCola label is still running. While domain1 is running, the run-time exclusion set of our policy says that domain2 cannot start because the label of domain1 includes the CHWALL type CocaCola and the label of domain2 includes the CHWALL type PepsiCo. The run-time exclusion rule of our policy enforces that PepsiCo and CocaCola cannot run at the same time on the same hypervisor platform. Once domain1 is stopped or saved, domain2 can start but domain1 can no longer start or be resumed. The ezPolicy tool,

when creating the Chinese Wall types for the workload labels, ensures that department workloads inherit the organization type (and with it any organization exclusions).

```
# xm list --label
Name           ID Mem(MiB) VCPUs State  Time(s)  Label
domain1        2    128      1 -b----- 6.9      CocaCola
Domain-0        0   1949      4 r----- 273.1    SystemManagement

# xm create domain2.xm
Using config file "domain2.xm".
Error: (1, 'Operation not permitted')

# xm destroy domain1
# xm create domain2.xm
Using config file "domain2.xm".
Started domain domain2

# xm list --label
Name           ID Mem(MiB) VCPUs State  Time(s)  Label
domain2        4    164      1 r----- 4.3      PepsiCo.Extranet
Domain-0        0   1949      4 r----- 298.4    SystemManagement

# xm create domain1.xm
Using config file "domain1.xm".
Error: (1, 'Operation not permitted')

# xm destroy domain2
# xm list
Name           ID Mem(MiB) VCPUs State  Time(s)
Domain-0        0   1949      4 r----- 391.2
```

You can verify that domains with Avis label can run together with domains labeled CocaCola, PepsiCo, or Hertz.

Test2: Resource access. In this test, we will re-label the swap file for domain1 with the Avis resource label. We expect that Domain1 will no longer start because it cannot access this resource. This test checks the sharing abilities of domains, which are defined by the Simple Type Enforcement Policy component.

```
# xm rmlabel res file:/home/xen/dom_fc5/fedora.swap
# xm addlabel Avis res file:/home/xen/dom_fc5/fedora.swap
# xm resources
file:/home/xen/dom_fc5/fedora.swap
  policy: example.chwall_ste.test-wld
  label:  Avis
file:/home/xen/dom_fc5/fedora.fc5.img
  policy: example.chwall_ste.test-wld
  label:  CocaCola

# xm create domain1.xm
Using config file "domain1.xm".
  file:/home/xen/dom_fc4/fedora.swap: DENIED
--> res:Avis (example.chwall_ste.test-wld)
--> dom:CocaCola (example.chwall_ste.test-wld)
Security configuration prevents domain from starting
```

Test 3: Communication. In this test we would verify that two domains with labels Hertz and Avis cannot exchange network packets by using the 'ping' connectivity test. It is also related to the STE policy.**Note:** sHype/Xen does control direct communication between domains. However, domains associated with different workloads can currently still communicate through the Domain0 virtual network. We are working on the sHype/ACM controls for local and remote network traffic through Domain0. Please monitor the xen-devel mailing list for updated information.

10.3 Xen Access Control Policy

This section describes the sHype/Xen access control policy in detail. It gives enough information to enable the reader to write custom access control policies and to use the available Xen policy tools. The policy language is expressive enough to specify most symmetric access relationships between domains and resources efficiently.

The Xen access control policy consists of two policy components. The first component, called Chinese Wall (CHWALL) policy, controls which domains can run simultaneously on the same virtualized platform. The second component, called Simple Type Enforcement (STE) policy, controls the sharing between running domains, i.e., communication or access to shared resources. The CHWALL and STE policy components can be configured to run alone, however in our examples we will assume that both policy components are configured together since they complement each other. The XML policy file includes all information needed by Xen to enforce the policies.

Figures 10.3 and 10.4 show a fully functional but very simple example policy for Xen. The policy can distinguish two workload types `CocaCola` and `PepsiCo` and defines the labels necessary to associate domains and resources with one of these workload types. The XML Policy consists of four parts:

1. policy header including the policy name
2. Simple Type Enforcement block
3. Chinese Wall Policy block
4. label definition block

10.3.1 Policy Header and Policy Name

Lines 1-2 (cf Figure 10.3) include the usual XML header. The security policy definition starts in Line 3 and refers to the policy schema. The XML-Schema definition for the Xen policy can be found in the file `/etc/xen/acm-security/policies/security-policy.xsd`. Examples for security policies can be found in the example subdirectory. The acm-security directory is only installed if ACM security is configured during installation (cf Section 10.2.1).

```

01 <?xml version="1.0" encoding="UTF-8"?>
02 <!-- Auto-generated by ezPolicy -->
03 <SecurityPolicyDefinition
    xmlns="http://www.ibm.com"
    xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
    xsi:schemaLocation=
        "http://www.ibm.com ../../security_policy.xsd ">
04     <PolicyHeader>
05         <PolicyName>example.chwall_ste.test</PolicyName>
06         <Date>Wed Jul 12 17:32:59 2006</Date>
07     </PolicyHeader>
08
09     <SimpleTypeEnforcement>
10         <SimpleTypeEnforcementTypes>
11             <Type>SystemManagement</Type>
12             <Type>PepsiCo</Type>
13             <Type>CocaCola</Type>
14         </SimpleTypeEnforcementTypes>
15     </SimpleTypeEnforcement>
16
17     <ChineseWall priority="PrimaryPolicyComponent">
18         <ChineseWallTypes>
19             <Type>SystemManagement</Type>
20             <Type>PepsiCo</Type>
21             <Type>CocaCola</Type>
22         </ChineseWallTypes>
23
24         <ConflictSets>
25             <Conflict name="RER1">
26                 <Type>CocaCola</Type>
27                 <Type>PepsiCo</Type>
28             </Conflict>
29         </ConflictSets>
30     </ChineseWall>
31

```

Figure 10.3: Example XML security policy file – Part I: Types and Rules Definition.

The `Policy Header` spans lines 4-7. It includes a date field and defines the policy name `example.chwall_ste.test`. It can also include optional fields that are not shown and are for future use (see schema definition).

The policy name serves two purposes: First, it provides a unique name for the security policy. This name is also exported by the Xen hypervisor to the Xen management tools in order to ensure that both enforce the same policy. We plan to extend the policy name with a digital fingerprint of the policy contents to better protect this correlation. Second, it implicitly points the xm tools to the location where the XML policy file is stored on the Xen system. Replacing the colons in the policy name by slashes yields the local path to the policy file starting from the global policy directory `/etc/xen/acm-security/policies`. The last part of the policy name is the prefix for the XML policy file name, completed by `-security_policy.xml`. Consequently, the policy with the name `example.chwall_ste.test` can be found in the XML policy file named `test-security_policy.xml` that is stored in the local directory `example/chwall_ste` under the global policy directory.

10.3.2 Simple Type Enforcement Policy Component

The Simple Type Enforcement (STE) policy controls which domains can communicate or share resources. This way, Xen can enforce confinement of workload types by confining the domains running those workload types. The mandatory access control framework enforces its policy when domains access intended ways of communication or cooperation (shared memory, events, shared resources such as block devices). It builds on top of the core hypervisor isolation, which restricts the ways of intercommunication to those intended means. STE does not protect or intend to protect from covert channels in the hypervisor or hardware; this is an orthogonal problem that can be mitigated by using the Run-time Exclusion rules described above or by fixing the problem in the core hypervisor.

Xen controls sharing between domains on the resource and domain level because this is the abstraction the hypervisor and its management understand naturally. While this is coarse-grained, it is also very reliable and robust and it requires minimal changes to implement mandatory access controls in the hypervisor. It enables platform- and operation system-independent policies as part of a layered security approach.

Lines 9-15 (cf Figure 10.3) define the Simple Type Enforcement policy component. Essentially, they define the workload type names `SystemManagement`, `PepsiCo`, and `CocaCola` that are available in the STE policy component. The policy rules are implicit: Xen permits a domain to communicate with another domain if and only if the labels of the domains share a common STE type. Xen permits a domain to access a resource if and only if the labels of the domain and the resource share a common STE workload type.

10.3.3 Chinese Wall Policy Component

The Chinese Wall security policy interpretation of sHype enables users to prevent certain workloads from running simultaneously on the same hypervisor platform. Run-time Exclusion rules (RER), also called Conflict Sets, define a set of workload types that are not permitted to run simultaneously. Of all the workloads specified in a Run-time Exclusion rule, at most one type can run on the same hypervisor platform at a time. Run-time Exclusion Rules implement a less rigorous variant of the original Chinese Wall security component. They do not implement the *-property of the policy, which would require to restrict also types that are not part of an exclusion rule once they are running together with a type in an exclusion rule (please refer to <http://www.gammasl.co.uk/topics/chinesewall.html> for more information on the original Chinese Wall policy).

Xen considers the `ChineseWallTypes` part of the label for the enforcement of the Run-time Exclusion rules. It is illegal to define labels including conflicting Chinese Wall types.

Lines 17-30 (cf Figure 10.3) define the Chinese Wall policy component. Lines 17-22 define the known Chinese Wall types, which coincide here with the STE types defined above. This usually holds if the criteria for sharing among domains and sharing of the hardware platform are the same. Lines 24-29 define one Run-time Exclusion rule:

```
<Conflict name="RER1">
  <Type>CocaCola</Type>
  <Type>PepsiCo</Type>
</Conflict>
```

Based on this rule, Xen enforces that only one of the types `CocaCola` or `PepsiCo` will run on a single hypervisor platform at a time. For example, once a domain assigned a `CocaCola` workload type is started, domains with the `PepsiCo` type will be denied to start. When the former domain stops and no other domains with the `CocaCola` type are running, then domains with the `PepsiCo` type can start.

Xen maintains reference counts on each running workload type to keep track of which workload types are running. Every time a domain starts or resumes, the reference count on those Chinese Wall types that are referenced in the domain's label are incremented. Every time a domain is destroyed or saved, the reference counts of its Chinese Wall types are decremented. `sHype` in Xen covers migration and live-migration, which is treated the same way as saving a domain on the source platform and resuming it on the destination platform.

Reasons why users would want to restrict which workloads or domains can share the system hardware include:

- Imperfect resource management or control might enable a rogue domain to starve another domain and the workload running in it.
- Redundant domains might run the same workload to increase availability; such domains should not run on the same hardware to avoid single points of failure.
- Imperfect Xen core domain isolation might enable two rogue domains running different workload types to use unintended and unknown ways (covert channels) to exchange some data. This way, they bypass the policed Xen access control mechanisms. Such imperfections cannot be completely eliminated and are a result of trade-offs between security and other design requirements. For a simple example of a covert channel see <http://www.multicians.org/timing-chn.html>. Such covert channels exist also between workloads running on different platforms if they are connected through networks. The Xen Chinese Wall policy provides an approximation of this imperfect “air-gap” between selected workload types.

10.3.4 Security Labels

To enable Xen to associate domains with workload types running in them, each domain is assigned a security label that includes the workload types of the domain.

```

32     <SecurityLabelTemplate>
33         <SubjectLabels bootstrap="SystemManagement">
34             <VirtualMachineLabel>
35                 <Name>SystemManagement</Name>
36                 <SimpleTypeEnforcementTypes>
37                     <Type>SystemManagement</Type>
38                     <Type>PepsiCo</Type>
39                     <Type>CocaCola</Type>
40                 </SimpleTypeEnforcementTypes>
41                 <ChineseWallTypes>
42                     <Type>SystemManagement</Type>
43                 </ChineseWallTypes>
44             </VirtualMachineLabel>
45
46             <VirtualMachineLabel>
47                 <Name>PepsiCo</Name>
48                 <SimpleTypeEnforcementTypes>
49                     <Type>PepsiCo</Type>
50                 </SimpleTypeEnforcementTypes>
51                 <ChineseWallTypes>
52                     <Type>PepsiCo</Type>
53                 </ChineseWallTypes>
54             </VirtualMachineLabel>
55
56             <VirtualMachineLabel>
57                 <Name>CocaCola</Name>
58                 <SimpleTypeEnforcementTypes>
59                     <Type>CocaCola</Type>
60                 </SimpleTypeEnforcementTypes>
61                 <ChineseWallTypes>
62                     <Type>CocaCola</Type>
63                 </ChineseWallTypes>
64             </VirtualMachineLabel>
65         </SubjectLabels>
66
67         <ObjectLabels>
68             <ResourceLabel>
69                 <Name>SystemManagement</Name>
70                 <SimpleTypeEnforcementTypes>
71                     <Type>SystemManagement</Type>
72                 </SimpleTypeEnforcementTypes>
73             </ResourceLabel>
74
75             <ResourceLabel>
76                 <Name>PepsiCo</Name>
77                 <SimpleTypeEnforcementTypes>
78                     <Type>PepsiCo</Type>
79                 </SimpleTypeEnforcementTypes>
80             </ResourceLabel>
81
82             <ResourceLabel>
83                 <Name>CocaCola</Name>
84                 <SimpleTypeEnforcementTypes>
85                     <Type>CocaCola</Type>
86                 </SimpleTypeEnforcementTypes>
87             </ResourceLabel>
88         </ObjectLabels>
89     </SecurityLabelTemplate>
90 </SecurityPolicyDefinition>

```

Figure 10.4: Example XML security policy file – Part II: Label Definition.

Lines 32-89 (cf Figure 10.4) define the `SecurityLabelTemplate`, which includes the labels that can be attached to domains and resources when this policy is active. The domain labels include Chinese Wall types while resource labels do not include Chinese Wall types. Lines 33-65 define the `SubjectLabels` that can be assigned to domains. For example, the virtual machine label `CocaCola` (cf lines 56-64 in Figure 10.4) associates the domain that carries it with the workload type `CocaCola`.

The `bootstrap` attribute names the label `SystemManagement`. Xen will assign this label to `Domain0` at boot time. All other domains are assigned labels according to their domain configuration file (see Section 10.2.4 for examples of how to label domains). Lines 67-88 define the `ObjectLabels`. Those labels can be assigned to resources when this policy is active.

In general, user domains should be assigned labels that have only a single `SimpleTypeEnforcement` workload type. This way, workloads remain confined even if user domains become rogue. Any domain that is assigned a label with multiple STE types must be trusted to keep information belonging to the different STE types separate (confined). For example, `Domain0` is assigned the bootstrap label `SystemsManagement`, which includes all existing STE types. Therefore, `Domain0` must take care not to enable unauthorized information flow (eg. through block devices or virtual networking) between domains or resources that are assigned different STE types.

Security administrators simply use the name of a label (specified in the `<Name>` field) to associate a label with a domain (cf. Section 10.2.4). The types inside the label are used by the Xen access control enforcement. While the name can be arbitrarily chosen (as long as it is unique), it is advisable to choose the label name in accordance to the security types included. While the XML representation in the above label seems unnecessary flexible, labels in general can consist of multiple types as we will see in the following example.

Assume that `PepsiCo` and `CocaCola` workloads use virtual disks that are provided by a virtual I/O domain hosting a physical storage device and carrying the following label:

```
<VirtualMachineLabel>
  <Name>VIO</Name>
  <SimpleTypeEnforcementTypes>
    <Type>CocaCola</Type>
    <Type>PepsiCo</Type>
  </SimpleTypeEnforcementTypes>
  <ChineseWallTypes>
    <Type>VIOServer</Type>
  </ChineseWallTypes>
</VirtualMachineLabel>
```

This Virtual I/O domain (VIO) exports its virtualized disks by communicating both to domains labeled with the `PepsiCo` label and domains labeled with the `CocaCola` label. This requires the VIO domain to carry both the STE types `CocaCola` and

PepsiCo. In this example, the confinement of CocaCola and PepsiCo workload depends on a VIO domain that must keep the data of those different workloads separate. The virtual disks are labeled as well (see Section 10.2.5 for labeling resources) and enforcement functions inside the VIO domain must ensure that the labels of the domain mounting a virtual disk and the virtual disk label share a common STE type. The VIO label carrying its own VIOServer CHWALL type introduces the flexibility to permit the trusted VIO server to run together with CocaCola or PepsiCo workloads.

Alternatively, a system that has two hard-drives does not need a VIO domain but can directly assign one hardware storage device to each of the workloads (if the platform offers an IO-MMU, cf Section 9.2. Sharing hardware through virtualization is a trade-off between the amount of trusted code (size of the trusted computing base) and the amount of acceptable over-provisioning. This holds both for peripherals and for system platforms.

10.3.5 Tools For Creating sHype/Xen Security Policies

To create a security policy for Xen, you can use one of the following tools:

- `ezPolicy` GUI tool – start writing policies
- `xensec_gen` tool – refine policies created with `ezPolicy`
- text or XML editor

We use the `ezPolicy` tool in Section 10.2.2 to quickly create a workload protection policy. If desired, the resulting XML policy file can be loaded into the `xensec_gen` tool to refine it. It can also be directly edited using an XML editor. Any XML policy file is verified against the security policy schema when it is translated (see Subsection 10.2.3).

10.4 Current Limitations

The sHype/ACM configuration for Xen is work in progress. There is ongoing work for protecting virtualized resources and planned and ongoing work for protecting access to remote resources and domains. The following sections describe limitations of some of the areas into which access control is being extended.

10.4.1 Network Traffic

Local and remote network traffic is currently not controlled. Solutions to add sHype/ACM policy enforcement to the virtual network exist but need to be discussed before they can become part of Xen. Subjecting external network traffic to the ACM security pol-

icy is work in progress. Manually setting up filters in domain 0 is required for now but does not scale well.

10.4.2 Resource Access and Usage Control

Enforcing the security policy across multiple hypervisor systems and on access to remote shared resources is work in progress. Extending access control to new types of resources is ongoing work (e.g. network storage).

On a single Xen system, information about the association of resources and security labels is stored in `/etc/xen/acm-security/policy/resource_labels`. This file relates a full resource path with a security label. This association is weak and will break if resources are moved or renamed without adapting the label file. Improving the protection of label-resource relationships is ongoing work.

Controlling resource usage and enforcing resource limits in general is ongoing work in the Xen community.

10.4.3 Domain Migration

Labels on domains are enforced during domain migration and the destination hypervisor will ensure that the domain label is valid and the domain is permitted to run (considering the Chinese Wall policy rules) before it accepts the migration. However, the network between the source and destination hypervisor as well as both hypervisors must be trusted. Architectures and prototypes exist that both protect the network connection and ensure that the hypervisors enforce access control consistently but patches are not yet available for the main stream.

10.4.4 Covert Channels

The sHype access control aims at system independent security policies. It builds on top of the core hypervisor isolation. Any covert channels that exist in the core hypervisor or in the hardware (e.g., shared processor cache) will be inherited. If those covert channels are not the result of trade-offs between security and other system properties, then they are most effectively minimized or eliminated where they are caused. sHype offers however some means to mitigate their impact (cf. run-time exclusion rules).

Part III

Reference

Chapter 11

Build and Boot Options

This chapter describes the build- and boot-time options which may be used to tailor your Xen system.

11.1 Top-level Configuration Options

Top-level configuration is achieved by editing one of two files: `Config.mk` and `Makefile`.

The former allows the overall build target architecture to be specified. You will typically not need to modify this unless you are cross-compiling or if you wish to build a PAE-enabled Xen system. Additional configuration options are documented in the `Config.mk` file.

The top-level `Makefile` is chiefly used to customize the set of kernels built. Look for the line:

```
KERNELS ?= linux-2.6-xen0 linux-2.6-xenU
```

Allowable options here are any kernels which have a corresponding build configuration file in the `buildconfigs/` directory.

11.2 Xen Build Options

Xen provides a number of build-time options which should be set as environment variables or passed on make's command-line.

verbose=y Enable debugging messages when Xen detects an unexpected condition. Also enables console output from all domains.

debug=y Enable debug assertions. Implies **verbose=y**. (Primarily useful for tracing bugs in Xen).

debugger=y Enable the in-Xen debugger. This can be used to debug Xen, guest OSes, and applications.

perf=y Enable performance counters for significant events within Xen. The counts can be reset or displayed on Xen's console via console control keys.

11.3 Xen Boot Options

These options are used to configure Xen's behaviour at runtime. They should be appended to Xen's command line, either manually or by editing `grub.conf`.

noreboot Don't reboot the machine automatically on errors. This is useful to catch debug output if you aren't catching console messages via the serial line.

nosmp Disable SMP support. This option is implied by 'ignorebiosables'.

watchdog Enable NMI watchdog which can report certain failures.

noirqbalance Disable software IRQ balancing and affinity. This can be used on systems such as Dell 1850/2850 that have workarounds in hardware for IRQ-routing issues.

badpage=<page number>,<page number>,... Specify a list of pages not to be allocated for use because they contain bad bytes. For example, if your memory tester says that byte 0x12345678 is bad, you would place 'badpage=0x12345' on Xen's command line.

com1=<baud>,<DPS>,<io_base>,<irq> com2=<baud>,<DPS>,<io_base>,<irq>

Xen supports up to two 16550-compatible serial ports. For example: 'com1=9600, 8n1, 0x408, 5' maps COM1 to a 9600-baud port, 8 data bits, no parity, 1 stop bit, I/O port base 0x408, IRQ 5. If some configuration options are standard (e.g., I/O base and IRQ), then only a prefix of the full configuration string need be specified. If the baud rate is pre-configured (e.g., by the bootloader) then you can specify 'auto' in place of a numeric baud rate.

console=<specifier list> Specify the destination for Xen console I/O. This is a comma-separated list of, for example:

vga Use VGA console (until domain 0 boots, unless **vga=keep** is specified).

com1 Use serial port com1.

com2H Use serial port com2. Transmitted chars will have the MSB set. Received chars must have MSB set.

com2L Use serial port com2. Transmitted chars will have the MSB cleared. Received chars must have MSB cleared.

The latter two examples allow a single port to be shared by two subsystems

(e.g. console and debugger). Sharing is controlled by MSB of each transmitted/received character. [NB. Default for this option is 'com1,vga']

vga=<options> This is a comma-separated list of options:

text-<mode> Select text-mode resolution, where mode is one of 80x25, 80x28, 80x30, 80x34, 80x43, 80x50, 80x60.

keep Keep the VGA console even after domain 0 boots.

sync_console Force synchronous console output. This is useful if your system fails unexpectedly before it has sent all available output to the console. In most cases Xen will automatically enter synchronous mode when an exceptional event occurs, but this option provides a manual fallback.

conswitch=<switch-char><auto-switch-char> Specify how to switch serial-console input between Xen and DOM0. The required sequence is CTRL-<switch-char> pressed three times. Specifying the backtick character disables switching. The <auto-switch-char> specifies whether Xen should auto-switch input to DOM0 when it boots — if it is 'x' then auto-switching is disabled. Any other value, or omitting the character, enables auto-switching. [NB. Default switch-char is 'a'.]

nmi=xxx Specify what to do with an NMI parity or I/O error.

'nmi=fatal': Xen prints a diagnostic and then hangs.

'nmi=dom0': Inform DOM0 of the NMI.

'nmi=ignore': Ignore the NMI.

mem=xxx Set the physical RAM address limit. Any RAM appearing beyond this physical address in the memory map will be ignored. This parameter may be specified with a B, K, M or G suffix, representing bytes, kilobytes, megabytes and gigabytes respectively. The default unit, if no suffix is specified, is kilobytes.

dom0_mem=xxx Set the amount of memory to be allocated to domain0. In Xen 3.x the parameter may be specified with a B, K, M or G suffix, representing bytes, kilobytes, megabytes and gigabytes respectively; if no suffix is specified, the parameter defaults to kilobytes. In previous versions of Xen, suffixes were not supported and the value is always interpreted as kilobytes.

dom0_vcpus_pin Pins domain 0 VCPUs on their respective physical CPUs (default=false).

tbuf_size=xxx Set the size of the per-cpu trace buffers, in pages (default 0).

sched=xxx Select the CPU scheduler Xen should use. The current possibilities are 'credit' (default), and 'sedf'.

apic_verbosity=debug,verbose Print more detailed information about local APIC and IOAPIC configuration.

lapic Force use of local APIC even when left disabled by uniprocessor BIOS.

nolapic Ignore local APIC in a uniprocessor system, even if enabled by the BIOS.

apic=bigsmpt, default, es7000, summit Specify NUMA platform. This can usually be probed automatically.

In addition, the following options may be specified on the Xen command line. Since domain 0 shares responsibility for booting the platform, Xen will automatically propagate these options to its command line. These options are taken from Linux's command-line syntax with unchanged semantics.

acpi=off, force, strict, ht, noirq, ... Modify how Xen (and domain 0) parses the BIOS ACPI tables.

acpi_skip_timer_override Instruct Xen (and domain 0) to ignore timer-interrupt override instructions specified by the BIOS ACPI tables.

noapic Instruct Xen (and domain 0) to ignore any IOAPICs that are present in the system, and instead continue to use the legacy PIC.

11.4 XenLinux Boot Options

In addition to the standard Linux kernel boot options, we support:

xencons=xxx Specify the device node to which the Xen virtual console driver is attached. The following options are supported:

‘xencons=off’: disable virtual console

‘xencons=tty’: attach console to /dev/tty1 (tty0 at boot-time)

‘xencons=ttyS’: attach console to /dev/ttyS0

The default is ttyS for dom0 and tty for all other domains.

Chapter 12

Further Support

If you have questions that are not answered by this manual, the sources of information listed below may be of interest to you. Note that bug reports, suggestions and contributions related to the software (or the documentation) should be sent to the Xen developers' mailing list (address below).

12.1 Other Documentation

For developers interested in porting operating systems to Xen, the *Xen Interface Manual* is distributed in the `docs/` directory of the Xen source distribution.

12.2 Online References

The official Xen web site can be found at:

`http://www.xensource.com`

This contains links to the latest versions of all online documentation, including the latest version of the FAQ.

Information regarding Xen is also available at the Xen Wiki at

`http://wiki.xensource.com/xenwiki/`

The Xen project uses Bugzilla as its bug tracking system. You'll find the Xen Bugzilla at `http://bugzilla.xensource.com/bugzilla/`.

12.3 Mailing Lists

There are several mailing lists that are used to discuss Xen related topics. The most widely relevant are listed below. An official page of mailing lists and subscription information can be found at

<http://lists.xensource.com/>

xen-devel@lists.xensource.com Used for development discussions and bug reports.

Subscribe at:

<http://lists.xensource.com/xen-devel>

xen-users@lists.xensource.com Used for installation and usage discussions and requests for help. Subscribe at:

<http://lists.xensource.com/xen-users>

xen-announce@lists.xensource.com Used for announcements only. Subscribe at:

<http://lists.xensource.com/xen-announce>

xen-changelog@lists.xensource.com Changelog feed from the unstable and 2.0 trees - developer oriented. Subscribe at:

<http://lists.xensource.com/xen-changelog>

Appendix A

Unmodified (VMX) guest domains in Xen with Intel® Virtualization Technology (VT)

Xen supports guest domains running unmodified Guest operating systems using Virtualization Technology (VT) available on recent Intel Processors. More information about the Intel Virtualization Technology implementing Virtual Machine Extensions (VMX) in the processor is available on the Intel website at

<http://www.intel.com/technology/computing/vptech>

A.1 Building Xen with VT support

The following packages need to be installed in order to build Xen with VT support. Some Linux distributions do not provide these packages by default.

Package	Description
dev86	<p>The dev86 package provides an assembler and linker for real mode 80x86 instructions. You need to have this package installed in order to build the BIOS code which runs in (virtual) real mode.</p> <p>If the dev86 package is not available on the x86_64 distribution, you can install the i386 version of it. The dev86 rpm package for various distributions can be found at http://www.rpmfind.net/linux/rpm2html/search.php?query=dev86&submit=Search</p>
LibVNCServer	<p>The unmodified guest's VGA display, keyboard, and mouse can be virtualized by the vncserver library. You can get the sources of libvncserver from http://sourceforge.net/projects/libvncserver. Build and install the sources on the build system to get the libvncserver library. There is a significant performance degradation in 0.8 version. The current sources in the CVS tree have fixed this degradation. So it is highly recommended to download the latest CVS sources and install them.</p>
SDL-devel, SDL	<p>Simple DirectMedia Layer (SDL) is another way of virtualizing the unmodified guest console. It provides an X window for the guest console. If the SDL and SDL-devel packages are not installed by default on the build system, they can be obtained from http://www.rpmfind.net/linux/rpm2html/search.php?query=SDL&submit=Search, http://www.rpmfind.net/linux/rpm2html/search.php?query=SDL-devel&submit=Search</p>

A.2 Configuration file for unmodified VMX guests

The Xen installation includes a sample configuration file, `/etc/xen/xmexample.vmx`. There are comments describing all the options. In addition to the common options that are the same as those for paravirtualized guest configurations, VMX guest configurations have the following settings:

Parameter	Description
kernel	The VMX firmware loader, <code>/usr/lib/xen/boot/vmxloader</code>
builder	The domain build function. The VMX domain uses the <code>vmx</code> builder.
acpi	Enable VMX guest ACPI, default=0 (disabled)
apic	Enable VMX guest APIC, default=0 (disabled)
paе	Enable VMX guest PAE, default=0 (disabled)
vif	Optionally defines MAC address and/or bridge for the network interfaces. Random MACs are assigned if not given. <code>type=ioemu</code> means <code>ioemu</code> is used to virtualize the VMX NIC. If no type is specified, <code>vbd</code> is used, as with paravirtualized guests.
disk	<p>Defines the disk devices you want the domain to have access to, and what you want them accessible as. If using a physical device as the VMX guest's disk, each disk entry is of the form</p> <pre>phy:UNAME,ioemu:DEV,MODE,</pre> <p>where UNAME is the device, DEV is the device name the domain will see, and MODE is <code>r</code> for read-only, <code>w</code> for read-write. <code>ioemu</code> means the disk will use <code>ioemu</code> to virtualize the VMX disk. If not adding <code>ioemu</code>, it uses <code>vbd</code> like paravirtualized guests.</p> <p>If using disk image file, its form should be like</p> <pre>file:FILEPATH,ioemu:DEV,MODE</pre> <p>If using more than one disk, there should be a comma between each disk entry. For example:</p> <pre>disk = ['file:/var/images/image1.img,ioemu:hda,w', 'file:/var/images/image2.img,ioemu:hdb,w']</pre>
cdrom	Disk image for CD-ROM. The default is <code>/dev/cdrom</code> for Domain0. Inside the VMX domain, the CD-ROM will available as device <code>/dev/hdc</code> . The entry can also point to an ISO file.
boot	<p>Boot from floppy (a), hard disk (c) or CD-ROM (d). For example, to boot from CD-ROM, the entry should be:</p> <pre>boot='d'</pre>
device_model	The device emulation tool for VMX guests. This parameter should not be changed.
sdl	Enable SDL library for graphics, default = 0 (disabled)
vnc	Enable VNC library for graphics, default = 1 (enabled)
vncviewer	<p>Enable spawning of the <code>vncviewer</code> (only valid when <code>vnc=1</code>), default = 1 (enabled)</p> <p>If <code>vnc=1</code> and <code>vncviewer=0</code>, user can use <code>vncviewer</code> to manually connect VMX from remote. For example:</p> <pre>vncviewer domain0_IP_address:VMX_domain_id</pre>
ne2000	Enable <code>ne2000</code> , default = 0 (disabled; use <code>pcnet</code>)
serial	Enable redirection of VMX serial output to <code>pty</code> device

usb	Enable USB support without defining a specific USB device. This option defaults to 0 (disabled) unless the option <code>usbdevice</code> is specified in which case this option then defaults to 1 (enabled).
usbdevice	<p>Enable USB support and also enable support for the given device. Devices that can be specified are <code>mouse</code> (a PS/2 style mouse), <code>tablet</code> (an absolute pointing device) and <code>host:id1:id2</code> (a physical USB device on the host machine whose ids are <code>id1</code> and <code>id2</code>). The advantage of <code>tablet</code> is that Windows guests will automatically recognize and support this device so specifying the config line</p> <pre>usbdevice='tablet'</pre> <p>will create a mouse that works transparently with Windows guests under VNC. Linux doesn't recognize the USB tablet yet so Linux guests under VNC will still need the Summagraphics emulation. Details about mouse emulation are provided in section A.4.3.</p>
localtime	Set the real time clock to local time [default=0, that is, set to UTC].
enable-audio	Enable audio support. This is under development.
full-screen	Start in full screen. This is under development.
nographic	Another way to redirect serial output. If enabled, no ' <code>sdl</code> ' or ' <code>vnc</code> ' can work. Not recommended.

A.3 Creating virtual disks from scratch

A.3.1 Using physical disks

If you are using a physical disk or physical disk partition, you need to install a Linux OS on the disk first. Then the boot loader should be installed in the correct place. For example `/dev/sda` for booting from the whole disk, or `/dev/sda1` for booting from partition 1.

A.3.2 Using disk image files

You need to create a large empty disk image file first; then, you need to install a Linux OS onto it. There are two methods you can choose. One is directly installing it using a VMX guest while booting from the OS installation CD-ROM. The other is copying an installed OS into it. The boot loader will also need to be installed.

To create the image file:

The image size should be big enough to accommodate the entire OS. This example assumes the size is 1G (which is probably too small for most OSes).

```
# dd if=/dev/zero of=hd.img bs=1M count=1 seek=1023
```

To directly install Linux OS into an image file using a VMX guest:

Install Xen and create VMX with the original image file with booting from CD-ROM. Then it is just like a normal Linux OS installation. The VMX configuration file should have these two entries before creating:

```
cdrom='/dev/cdrom' boot='d'
```

If this method does not succeed, you can choose the following method of copying an installed Linux OS into an image file.

To copy a installed OS into an image file:

Directly installing is an easier way to make partitions and install an OS in a disk image file. But if you want to create a specific OS in your disk image, then you will most likely want to use this method.

1. Install a normal Linux OS on the host machine

You can choose any way to install Linux, such as using yum to install Red Hat Linux or YAST to install Novell SuSE Linux. The rest of this example assumes the Linux OS is installed in `/var/guestos/`.

2. Make the partition table

The image file will be treated as hard disk, so you should make the partition table in the image file. For example:

```
# losetup /dev/loop0 hd.img
# fdisk -b 512 -C 4096 -H 16 -S 32 /dev/loop0
press 'n' to add new partition
press 'p' to choose primary partition
press '1' to set partition number
press "Enter" keys to choose default value of "First Cylinder" parameter.
press "Enter" keys to choose default value of "Last Cylinder" parameter.
press 'w' to write partition table and exit
# losetup -d /dev/loop0
```

3. Make the file system and install grub

```
# ln -s /dev/loop0 /dev/loop
# losetup /dev/loop0 hd.img
# losetup -o 16384 /dev/loop1 hd.img
# mkfs.ext3 /dev/loop1
# mount /dev/loop1 /mnt
# mkdir -p /mnt/boot/grub
# cp /boot/grub/stage* /boot/grub/e2fs_stage1_5 /mnt/boot/grub
```

```
# umount /mnt
# grub
grub> device (hd0) /dev/loop
grub> root (hd0,0)
grub> setup (hd0)
grub> quit
# rm /dev/loop
# losetup -d /dev/loop0
# losetup -d /dev/loop1
```

The `losetup` option `-o 16384` skips the partition table in the image file. It is the number of sectors times 512. We need `/dev/loop` because `grub` is expecting a disk device *name*, where *name* represents the entire disk and *name1* represents the first partition.

4. Copy the OS files to the image

If you have Xen installed, you can easily use `lomount` instead of `losetup` and `mount` when coping files to some partitions. `lomount` just needs the partition information.

```
# lomount -t ext3 -diskimage hd.img -partition 1 /mnt/guest
# cp -ax /var/guestos/{root,dev,var,etc,usr,bin,sbin,lib} /mnt/guest
# mkdir /mnt/guest/{proc,sys,home,tmp}
```

5. Edit the `/etc/fstab` of the guest image

The `fstab` should look like this:

```
# vim /mnt/guest/etc/fstab
/dev/hda1 / ext3 defaults 1 1
none /dev/pts devpts gid=5,mode=620 0 0
none /dev/shm tmpfs defaults 0 0
none /proc proc defaults 0 0
none /sys sysfs defaults 0 0
```

6. umount the image file

```
# umount /mnt/guest
```

Now, the guest OS image `hd.img` is ready. You can also reference <http://free.oszoo.org> for quickstart images. But make sure to install the boot loader.

A.3.3 Install Windows into an Image File using a VMX guest

In order to install a Windows OS, you should keep `acpi=0` in your VMX configuration file.

A.4 VMX Guests

A.4.1 Editing the Xen VMX config file

Make a copy of the example VMX configuration file `/etc/xen/xmeaxmple.vmx` and edit the line that reads

```
disk = [ 'file:/var/images/guest.img,ioemu:hda,w' ]
```

replacing *guest.img* with the name of the guest OS image file you just made.

A.4.2 Creating VMX guests

Simply follow the usual method of creating the guest, using the `-f` parameter and providing the filename of your VMX configuration file:

```
# xend start
# xm create /etc/xen/vmxguest.vmx
```

In the default configuration, VNC is on and SDL is off. Therefore VNC windows will open when VMX guests are created. If you want to use SDL to create VMX guests, set `sdl=1` in your VMX configuration file. You can also turn off VNC by setting `vnc=0`.

A.4.3 Mouse issues, especially under VNC

Mouse handling when using VNC is a little problematic. The problem is that the VNC viewer provides a virtual pointer which is located at an absolute location in the VNC window and only absolute coordinates are provided. The VMX device model converts these absolute mouse coordinates into the relative motion deltas that are expected by the PS/2 mouse driver running in the guest. Unfortunately, it is impossible to keep these generated mouse deltas accurate enough for the guest cursor to exactly match the VNC pointer. This can lead to situations where the guest's cursor is in the center of the screen and there's no way to move that cursor to the left (it can happen that the VNC pointer is at the left edge of the screen and, therefore, there are no longer any left mouse deltas that can be provided by the device model emulation code.)

To deal with these mouse issues there are 4 different mouse emulations available from the VMX device model:

PS/2 mouse over the PS/2 port. This is the default mouse that works perfectly well under SDL. Under VNC the guest cursor will get out of sync with the VNC pointer. When this happens you can re-synchronize the guest cursor to the VNC pointer by holding down the **left-ctl** and **left-alt** keys together. While these keys are down VNC pointer motions will not be reported to the guest so that the

VNC pointer can be moved to a place where it is possible to move the guest cursor again.

Summagraphics mouse over the serial port. The device model also provides emulation for a Summagraphics tablet, an absolute pointer device. This emulation is provided over the second serial port, **/dev/ttyS1** for Linux guests and **COM2** for Windows guests. Unfortunately, neither Linux nor Windows provides default support for the Summagraphics tablet so the guest will have to be manually configured for this mouse.

Linux configuration.

First, configure the GPM service to use the Summagraphics tablet. This can vary between distributions but, typically, all that needs to be done is modify the file `/etc/sysconfig/mouse` to contain the lines:

```
MOUSETYPE="summa"  
XMOUSETYPE="SUMMA"  
DEVICE=/dev/ttyS1
```

and then restart the GPM daemon.

Next, modify the X11 config `/etc/X11/xorg.conf` to support the Summagraphics tablet by replacing the input device stanza with the following:

```
Section "InputDevice"  
    Identifier "Mouse0"  
    Driver "summa"  
    Option "Device" "/dev/ttyS1"  
    Option "InputFashion" "Tablet"  
    Option "Mode" "Absolute"  
    Option "Name" "EasyPen"  
    Option "Compatible" "True"  
    Option "Protocol" "Auto"  
    Option "SendCoreEvents" "on"  
    Option "Vendor" "GENIUS"  
EndSection
```

Restart X and the X cursor should now properly track the VNC pointer.

Windows configuration.

Get the file <http://www.cad-plan.de/files/download/tw2k.exe> and execute that file on the guest, answering the questions as follows:

1. When the program asks for **model**, scroll down and select **SummaSketch (MM Compatible)**.
2. When the program asks for **COM Port** specify **com2**.
3. When the program asks for a **Cursor Type** specify **4 button cursor/puck**.

4. The guest system will then reboot and, when it comes back up, the guest cursor will now properly track the VNC pointer.

PS/2 mouse over USB port. This is just the same PS/2 emulation except it is provided over a USB port. This emulation is enabled by the configuration flag:

```
usbdevice='mouse'
```

USB tablet over USB port. The USB tablet is an absolute pointing device that has the advantage that it is automatically supported under Windows guests, although Linux guests still require some manual configuration. This mouse emulation is enabled by the configuration flag:

```
usbdevice='tablet'
```

Linux configuration.

Unfortunately, there is no GPM support for the USB tablet at this point in time. If you intend to use a GPM pointing device under VNC you should configure the guest for Summagraphics emulation.

Support for X11 is available by following the instructions at <http://stz-softwaretechnik.com/~ke/touchscreen/evtouch.html> with one minor change. The `xorg.conf` given in those instructions uses the wrong values for the X & Y minimums and maximums, use the following config stanza instead:

```
Section "InputDevice"
    Identifier      "Tablet"
    Driver          "evtouch"
    Option          "Device"    "/dev/input/event2"
    Option          "DeviceName" "touchscreen"
    Option          "MinX"      "0"
    Option          "MinY"      "0"
    Option          "MaxX"      "32256"
    Option          "MaxY"      "32256"
    Option          "ReportingMode" "Raw"
    Option          "Emulate3Buttons"
    Option          "Emulate3Timeout" "50"
    Option          "SendCoreEvents" "On"
EndSection
```

Windows configuration.

Just enabling the USB tablet in the guest's configuration file is sufficient, Windows will automatically recognize and configure device drivers for this pointing device.

A.4.4 USB Support

There is support for an emulated USB mouse, an emulated USB tablet and physical low speed USB devices (support for high speed USB 2.0 devices is still under development).

USB PS/2 style mouse. Details on the USB mouse emulation are given in sections **A.2** and **A.4.3**. Enabling USB PS/2 style mouse emulation is just a matter of adding the line

```
usbdevice='mouse'
```

to the configuration file.

USB tablet. Details on the USB tablet emulation are given in sections **A.2** and **A.4.3**. Enabling USB tablet emulation is just a matter of adding the line

```
usbdevice='tablet'
```

to the configuration file.

USB physical devices. Access to a physical (low speed) USB device is enabled by adding a line of the form

```
usbdevice='host:vid:pid'
```

into the the configuration file.¹ **vid** and **pid** are a product id and vendor id that uniquely identify the USB device. These ids can be identified in two ways:

1. Through the control window. As described in section **A.4.6** the control window is activated by pressing **ctl-alt-2** in the guest VGA window. As long as USB support is enabled in the guest by including the config file line

```
usb=1
```

then executing the command

```
info usbhost
```

in the control window will display a list of all usb devices and their ids. For example, this output:

```
Device 1.3, speed 1.5 Mb/s
Class 00: USB device 04b3:310b
```

was created from a USB mouse with vendor id **04b3** and product id **310b**. This device could be made available to the VMX guest by including the config file entry

```
usbdevice='host:04be:310b'
```

¹There is an alternate way of specifying a USB device that uses the syntax **host:bus.addr** but this syntax suffers from a major problem that makes it effectively useless. The problem is that the **addr** portion of this address changes every time the USB device is plugged into the system. For this reason this addressing scheme is not recommended and will not be documented further.

It is also possible to enable access to a USB device dynamically through the control window. The control window command

```
usb_add host:vid:pid
```

will also allow access to a USB device with vendor id **vid** and product id **pid**.

2. Through the `/proc` file system. The contents of the pseudo file `/proc/bus/usb/devices` can also be used to identify vendor and product ids. Looking at this file, the line starting with **P:** has a field **Vendor** giving the vendor id and another field **ProdID** giving the product id. The contents of `/proc/bus/usb/devices` for the example mouse is as follows:

```
T: Bus=01 Lev=01 Prnt=01 Port=01 Cnt=02 Dev#= 3 Spd=1.5 MxCh= 0
D: Ver= 2.00 Cls=00(>ifc ) Sub=00 Prot=00 MxPS= 8 #Cfgs= 1
P: Vendor=04b3 ProdID=310b Rev= 1.60
C:* #Ifs= 1 Cfg#= 1 Atr=a0 MxPwr=100mA
I: If#= 0 Alt= 0 #EPs= 1 Cls=03(HID ) Sub=01 Prot=02 Driver=(none)
E: Ad=81(I) Atr=03(Int.) MxPS= 4 Iv1=10ms
```

Note that the **P:** line correctly identifies the vendor id and product id for this mouse as **04b3:310b**.

There is one other issue to be aware of when accessing a physical USB device from the guest. The Dom0 kernel must not have a device driver loaded for the device that the guest wishes to access. This means that the Dom0 kernel must not have that device driver compiled into the kernel or, if using modules, that driver module must not be loaded. Note that this is the device specific USB driver that must not be loaded, either the **UHCI** or **OHCI** USB controller driver must still be loaded.

Going back to the USB mouse as an example, if **lsmod** gives the output:

Module	Size	Used by
usbmouse	4128	0
usbhid	28996	0
uhci_hcd	35409	0

then the USB mouse is being used by the Dom0 kernel and is not available to the guest. Executing the command **rmmod usbhid**² will remove the USB mouse driver from the Dom0 kernel and the mouse will now be accessible by the VMX guest.

Be aware the the Linux USB hotplug system will reload the drivers if a USB device is removed and plugged back in. This means that just unloading the driver module might not be sufficient if the USB device is removed and added back. A more reliable technique is to first **rmmod** the driver and then rename

²Turns out the **usbhid** driver is the significant one for the USB mouse, the presence or absence of the module **usbmouse** has no effect on whether or not the guest can see a USB mouse.

the driver file in the `/lib/modules` directory, just to make sure it doesn't get reloaded.

A.4.5 Destroy VMX guests

VMX guests can be destroyed in the same way as can paravirtualized guests. We recommend that you type the command

```
poweroff
```

in the VMX guest's console first to prevent data loss. Then execute the command

```
xm destroy vmx_guest_id
```

at the Domain0 console.

A.4.6 VMX window (X or VNC) Hot Key

If you are running in the X environment after creating a VMX guest, an X window is created. There are several hot keys for control of the VMX guest that can be used in the window.

Ctrl+Alt+2 switches from guest VGA window to the control window. Typing `help` shows the control commands help. For example, 'q' is the command to destroy the VMX guest.

Ctrl+Alt+1 switches back to VMX guest's VGA.

Ctrl+Alt+3 switches to serial port output. It captures serial output from the VMX guest. It works only if the VMX guest was configured to use the serial port.

A.4.7 Save/Restore and Migration

VMX guests currently cannot be saved and restored, nor migrated. These features are currently under active development.

Appendix B

Vnets - Domain Virtual Networking

Xen optionally supports virtual networking for domains using *vnets*. These emulate private LANs that domains can use. Domains on the same vnet can be hosted on the same machine or on separate machines, and the vnets remain connected if domains are migrated. Ethernet traffic on a vnet is tunneled inside IP packets on the physical network. A vnet is a virtual network and addressing within it need have no relation to addressing on the underlying physical network. Separate vnets, or vnets and the physical network, can be connected using domains with more than one network interface and enabling IP forwarding or bridging in the usual way.

Vnet support is included in `xm` and `xend`:

```
# xm vnet-create <config>
```

creates a vnet using the configuration in the file `<config>`. When a vnet is created its configuration is stored by `xend` and the vnet persists until it is deleted using

```
# xm vnet-delete <vnetid>
```

The vnets `xend` knows about are listed by

```
# xm vnet-list
```

More vnet management commands are available using the `vn` tool included in the vnet distribution.

The format of a vnet configuration file is

```
(vnet (id          <vnetid>)
      (bridge      <bridge>)
      (vnetif      <vnet interface>)
      (security    <level>))
```

White space is not significant. The parameters are:

- `<vnetid>`: vnet id, the 128-bit vnet identifier. This can be given as 8 4-digit hex numbers separated by colons, or in short form as a single 4-digit hex number. The short form is the same as the long form with the first 7 fields zero. Vnet ids must be non-zero and id 1 is reserved.
- `<bridge>`: the name of a bridge interface to create for the vnet. Domains are connected to the vnet by connecting their virtual interfaces to the bridge. Bridge names are limited to 14 characters by the kernel.
- `<vnetif>`: the name of the virtual interface onto the vnet (optional). The interface encapsulates and decapsulates vnet traffic for the network and is attached to the vnet bridge. Interface names are limited to 14 characters by the kernel.
- `<level>`: security level for the vnet (optional). The level may be one of
 - `none`: no security (default). Vnet traffic is in clear on the network.
 - `auth`: authentication. Vnet traffic is authenticated using IPSEC ESP with hmac96.
 - `conf`: confidentiality. Vnet traffic is authenticated and encrypted using IPSEC ESP with hmac96 and AES-128.

Authentication and confidentiality are experimental and use hard-wired keys at present.

When a vnet is created its configuration is stored by xend and the vnet persists until it is deleted using `xm vnet-delete <vnetid>`. The interfaces and bridges used by vnets are visible in the output of `ifconfig` and `brctl show`.

B.1 Example

If the file `vnet97.sxp` contains

```
(vnet (id 97) (bridge vnet97) (vnetif vnif97)
      (security none))
```

Then `xm vnet-create vnet97.sxp` will define a vnet with id 97 and no security. The bridge for the vnet is called `vnet97` and the virtual interface for it is `vnif97`. To add an interface on a domain to this vnet set its bridge to `vnet97` in its configuration. In Python:

```
vif="bridge=vnet97"
```

In `sxp`:

```
(dev (vif (mac aa:00:00:01:02:03) (bridge vnet97)))
```

Once the domain is started you should see its interface in the output of `brctl show` under the ports for `vnet97`.

To get best performance it is a good idea to reduce the MTU of a domain's interface onto a vnet to 1400. For example using `ifconfig eth0 mtu 1400` or putting `MTU=1400` in `ifcfg-eth0`. You may also have to change or remove cached config files for `eth0` under `/etc/sysconfig/networking`. Vnets work anyway, but performance can be reduced by IP fragmentation caused by the vnet encapsulation exceeding the hardware MTU.

B.2 Installing vnet support

Vnets are implemented using a kernel module, which needs to be loaded before they can be used. You can either do this manually before starting `xend`, using the command `vn insmod`, or configure `xend` to use the `network-vnet` script in the `xend` configuration file `/etc/xend/xend-config.sxp`:

```
(network-script          network-vnet)
```

This script `insmod`s the module and calls the `network-bridge` script.

The vnet code is not compiled and installed by default. To compile the code and install on the current system use `make install` in the root of the vnet source tree, `tools/vnet`. It is also possible to install to an installation directory using `make dist`. See the `Makefile` in the source for details.

The vnet module creates vnet interfaces `vnif0002`, `vnif0003` and `vnif0004` by default. You can test that vnets are working by configuring IP addresses on these interfaces and trying to ping them across the network. For example, using machines `hostA` and `hostB`:

```
hostA# ifconfig vnif0004 10.0.0.100 up
hostB# ifconfig vnif0004 10.0.0.101 up
hostB# ping 10.0.0.100
```

The vnet implementation uses IP multicast to discover vnet interfaces, so all machines hosting vnets must be reachable by multicast. Network switches are often configured not to forward multicast packets, so this often means that all machines using a vnet must be on the same LAN segment, unless you configure vnet forwarding.

You can test multicast coverage by pinging the vnet multicast address:

```
# ping -b 224.10.0.1
```

You should see replies from all machines with the vnet module running. You can see if vnet packets are being sent or received by dumping traffic on the vnet UDP port:

```
# tcpdump udp port 1798
```

If multicast is not being forwarded between machines you can configure multicast forwarding using `vn`. Suppose we have machines `hostA` on `10.10.0.100` and `hostB` on

10.11.0.100 and that multicast is not forwarded between them. We use vn to configure each machine to forward to the other:

```
hostA# vn peer-add hostB
```

```
hostB# vn peer-add hostA
```

Multicast forwarding needs to be used carefully - you must avoid creating forwarding loops. Typically only one machine on a subnet needs to be configured to forward, as it will forward multicasts received from other machines on the subnet.

Appendix C

Glossary of Terms

Domain A domain is the execution context that contains a running **virtual machine**.

The relationship between virtual machines and domains on Xen is similar to that between programs and processes in an operating system: a virtual machine is a persistent entity that resides on disk (somewhat like a program). When it is loaded for execution, it runs in a domain. Each domain has a **domain ID**.

Domain 0 The first domain to be started on a Xen machine. Domain 0 is responsible for managing the system.

Domain ID A unique identifier for a **domain**, analogous to a process ID in an operating system.

Full virtualization An approach to virtualization which requires no modifications to the hosted operating system, providing the illusion of a complete system of real hardware devices.

Hypervisor An alternative term for **VMM**, used because it means ‘beyond supervisor’, since it is responsible for managing multiple ‘supervisor’ kernels.

Live migration A technique for moving a running virtual machine to another physical host, without stopping it or the services running on it.

Paravirtualization An approach to virtualization which requires modifications to the operating system in order to run in a virtual machine. Xen uses paravirtualization but preserves binary compatibility for user space applications.

Shadow pagetables A technique for hiding the layout of machine memory from a virtual machine’s operating system. Used in some **VMMs** to provide the illusion of contiguous physical memory, in Xen this is used during **live migration**.

Virtual Block Device Persistent storage available to a virtual machine, providing the abstraction of an actual block storage device. **VBDs** may be actual block devices, filesystem images, or remote/network storage.

Virtual Machine The environment in which a hosted operating system runs, provid-

ing the abstraction of a dedicated machine. A virtual machine may be identical to the underlying hardware (as in **full virtualization**, or it may differ, as in **paravirtualization**).

VMM Virtual Machine Monitor - the software that allows multiple virtual machines to be multiplexed on a single physical machine.

Xen Xen is a paravirtualizing virtual machine monitor, developed primarily by the Systems Research Group at the University of Cambridge Computer Laboratory.

XenLinux A name for the port of the Linux kernel that runs on Xen.