# The Linux 2.4 SCSI subsystem HOWTO

**Douglas Gilbert** `<dgilbert at interlog dot com>`

# The Linux 2.4 SCSI subsystem HOWTO

by Douglas Gilbert

Publication date 2003-08-24
Copyright © 2001, 2002, 2003, 2004 Douglas Gilbert

**Abstract**

This document describes the SCSI subsystem as the Linux kernel enters the 2.4 production series. An external view of the SCSI subsystem is the main theme. Material is included to help the system administration of the Linux SCSI subsystem. There are also brief descriptions of ioctl()s and interfaces that may be relevant to those writing applications that use this subsystem.

# Table of Contents

# Chapter 1. Introduction

This document describes the SCSI subsystem as the Linux kernel enters the 2.4 production series.

An external view of the SCSI subsystem is the main theme. Material is included to help the system administration of the Linux SCSI subsystem. There are also brief descriptions of ioctl()s and interfaces that may be relevant to those writing applications that use this subsystem. However internal data structures and design issues are not addressed [see reference W2]. To unclutter the presentation, compile options and system calls (including ioctl()s) have been placed in Appendix E, *Compile options and System calls including ioctls*. Although not strictly part of the SCSI subsystem, there is also a description of raw devices in Chapter 11, *Raw devices*.

For those who have no interest in the SCSI subsystem and just want to get their ATAPI cd writer going, see the section called "ATAPI cdroms". It may also be useful to browse Chapter 2, *Architectural Overview*.

This document follows on from one written five years ago by Drew Eckhardt called the SCSI-HOWTO [see reference W7]. That document described the SCSI subsystem in Linux kernel 1.2 and 1.3 series. It is still available from the Linux Documentation Project [LDP, see reference W8] in its "unmaintained" section. Both documents have roughly similar structures although Drew's document has a lot of information on the adapter drivers.

This document can be found in electronic form at `www.tldp.org/HOWTO/SCSI-2.4-HOWTO` [http://www.tldp.org/HOWTO/SCSI-2.4-HOWTO]. The home site and perhaps the most up to date version of this document can be found at `www.torque.net/scsi/SCSI-2.4-HOWTO` [http://www.torque.net/scsi/SCSI-2.4-HOWTO] (this is the multi-page html version). At that location this document is rendered in txt, pdf, ps, a single (long) page of html as well as multi-page html. For example, a pdf version is at `www.torque.net/scsi/SCSI-2.4-HOWTO.pdf` [http://www.torque.net/scsi/SCSI-2.4-HOWTO.pdf]).

This document was last altered on 24th August 2004.

# Chapter 2. Architectural Overview

The SCSI subsystem has a 3 level architecture with the "upper" level being closest to the user/kernel interface while the "lower" level is closest to the hardware. The upper level drivers are commonly known by a terse two letter abbreviation (e.g. "sd" for SCSI disk driver). The names of the corresponding module drivers which, for historical reasons, sometimes differ from the built in driver names are shown in braces in the following diagram.

The 3 level driver architecture of the SCSI subsystem.

The upper level supports the user-kernel interface. In the case of sd and sr this is a block device interface while for st and sg this is a character device interface. Any operation using the SCSI subsystem (e.g. reading a sector from a disk) involves one driver at each of the 3 levels (e.g. sd, SCSI mid level and aic7xxx drivers).

As can be seen from the diagram, the SCSI mid level is common to all operations. The SCSI mid level defines internal interfaces and provides common services to the upper and lower level drivers. Ioctls provided by the mid level are available to the file descriptors belonging to any of the 4 upper level drivers.

The most common operation on a block device is to "mount" a file system. For a sd device typically a partition is mounted (e.g. **mount -t ext2 /dev/sda6 /home**). For a sr device usually the whole device is mounted (e.g. **mount -t iso9660 /dev/sr0 /mnt/cdrom**). The **dd** command can be used to read or write from block devices. In this case the block size argument ("bs") needs to be set to the block size of the device (e.g. 512 bytes for most disks) or an integral multiple of that device block size (e.g. 8192 bytes). A recent addition to the block subsystem allows a device (or partition) to be mounted more than once, at different mount points.

Sd is a member of the generic disk family, as is the hd device from the IDE subsystem. Apart from mounting sd devices, the **fdisk** command is available to view or modify a disk's partition table. Although the **hdparm** command is primarily intended for ATA disks (also known as IDE or EIDE disks), some options work on SCSI disks.

Sr is a member of the CD-ROM subsystem. Apart from mounting file systems (e.g. iso9660), audio CDs can also be read. The latter action does *not* involve mounting a file system but typically by invoking some ioctls. General purpose Linux commands such as **dd** cannot be used on audio CDs.

St is a char device for reading and writing tapes. Typically the **mt** command is used to perform data transfers and other control functions.

Sg is a SCSI command pass through device that uses a char device interface. General purpose Linux commands should *not* be used on sg devices. Applications such as SANE (for scanners), **cdrecord** and **cdrdao** (for cd writers) and **cdparanoia** (for reading audio CDs digitally) use sg.

# Chapter 3. Names and Addresses

This section covers the various naming schemes that exist in Linux and the SCSI worlds and how they interact.

## SCSI Addressing

Linux has a four level hierarchical addressing scheme for SCSI devices:

- SCSI adapter number [host]

- channel number [bus]

- id number [target]

- lun [lun]

"Lun" is the common SCSI abbreviation of Logical Unit Number. The terms in brackets are the name conventions used by device pseudo file system (devfs). "Bus" is used in preference to "channel" in the description below.

The SCSI adapter number is typically an arbitrary numbering of the adapter cards on the internal IO buses (e.g. PCI, PCMCIA, ISA etc) of the computer. Such adapters are sometimes termed as HBAs (host bus adapters). SCSI adapter numbers are issued by the kernel in ascending order starting with 0.

Each HBA may control one or more SCSI buses. The various types of SCSI buses are listed in Appendix A, *Common bus types (SCSI and other)*.

Each SCSI bus can have multiple SCSI devices connected to it. In SCSI parlance the HBA is called the "initiator" and takes up one SCSI id number (typically 7). The initiator [1] talks to targets which are commonly known as SCSI devices (e.g. disks). On SCSI parallel buses the number of ids is related to the width. 8 bit buses (sometimes called "narrow") can have 8 SCSI ids of which 1 is taken by the HBA leaving 7 for SCSI devices. Wide SCSI buses are 16 bits wide and can have a maximum of 15 SCSI devices (targets) attached. The SCSI 3 draft standard allows a large number of ids to be present on a SCSI bus.

Each SCSI device can contain multiple Logical Unit Numbers (LUNs). These are typically used by sophisticated tape and cdrom units that support multiple media.

So Linux's flavour of SCSI addressing is a four level hierarchy:

```
<scsi(_adapter_number), channel, id, lun>
```

Using the naming conventions of devfs this becomes:

```
<host, bus, target, lun>
```

## Device Names

A device name can be thought of as a gateway to a kernel driver that controls a device rather than the device itself. Hence there can be multiple device names some of which may offer slightly different characteristics, all mapping to the same actual device.

---

[1] SCSI standards allow for multiple initiators to be present on a single bus. This is not well supported in Linux although there are patches around that improve this situation.

The device names of the various SCSI devices are found within the /dev directory. Traditionally in Linux, SCSI devices have been identified by their major and minor device number rather than their SCSI bus addresses (e.g. SCSI target id and LUN). The device pseudo file system (devfs) moves away from the major and minor device number scheme and for the SCSI subsystem uses device names based on the SCSI bus addresses [discussed later in the section called "Device Names in devfs" and see reference: W5]. Alternatively, there is a utility called **scsidev** which addresses this issue within the scope of the Linux SCSI subsystem and thus does not have the same system wide impact as devfs. Scsidev is discussed later in the section called "Device Names in scsidev" and ref: W6.

Eight block major numbers are reserved for SCSI disks: 8, 65, 66, 67, 68, 69, 70 and 71. Each major can accommodate 256 minor numbers which, in the case of SCSI disks, are subdivided as follows:

```
[b,8,0]     /dev/sda
[b,8,1]     /dev/sda1
....
[b,8,15]    /dev/sda15
[b,8,16]    /dev/sdb
[b,8,17]    /dev/sdb1
....
[b,8,255]   /dev/sdp15
```

The disk device names without a trailing digit refer to the whole disk (e.g. /dev/sda) while those with a trailing digit refer to one of the 15 allowable partitions [2] within that disk.

The remaining 7 SCSI disk block major numbers follow a similar pattern:

```
[b,65,0]    /dev/sdq
[b,65,1]    /dev/sdq1
....
[b,65,159]  /dev/sdz15
[b,65,160]  /dev/sdaa
[b,65,161]  /dev/sdaa1
....
[b,65,255]  /dev/sdaf15
[b,66,0]    /dev/sdag
[b,66,1]    /dev/sdag1
....
[b,66,255]  /dev/sdav15
....
[b,71,255]  /dev/sddx15
```

So there are 128 possible disks (i.e. /dev/sda to /dev/sddx) each having up to 15 partitions. By way of contrast, the IDE subsystem allows 20 disks (10 controllers each with 1 master and 1 slave) which can have up to 63 partitions each.

SCSI CD-ROM devices are allocated the block major number of 11. Traditionally sr has been the device name but scd probably is more recognizable and is favoured by several recent distributions. 256 SCSI CD-ROM devices are allowed:

```
[b,11,0]    /dev/scd0              [or /dev/sr0]
```

---

[2] If 15 partitions is too limiting then the Logical Volume Manager (LVM) might be considered. See /usr/src/linux/Documentation/LVM-HOWTO . LVM will also allow a logical partition to span multiple block devices.

```
[b,11,255] /dev/scd255          [or /dev/sr255]
```

SCSI tape devices are allocated the char major number of 9. Up to 32 tape devices are supported each of which can be accessed in one of four modes (0, 1, 2 and 3), with or without rewind. The devices are allocated as follows:

```
[c,9,0]    /dev/st0     [tape 0, mode 0, rewind]
[c,9,1]    /dev/st1     [tape 1, mode 0, rewind]
....
[c,9,31]   /dev/st31    [tape 31, mode 0, rewind]
[c,9,32]   /dev/st0l    [tape 0, mode 1, rewind]
....
[c,9,63]   /dev/st31l   [tape 31, mode 1, rewind]
[c,9,64]   /dev/st0m    [tape 0, mode 2, rewind]
....
[c,9,96]   /dev/st0a    [tape 0, mode 3, rewind]
....
[c,9,127]  /dev/st31a   [tape 31, mode 3, rewind]
[c,9,128]  /dev/nst0    [tape 0, mode 0, no rewind]
....
[c,9,160]  /dev/nst0l   [tape 0, mode 1, no rewind]
....
[c,9,192]  /dev/nst0m   [tape 0, mode 2, no rewind]
....
[c,9,224]  /dev/nst0a   [tape 0, mode 3, no rewind]
....
[c,9,255]  /dev/nst31a  [tape 31, mode 3, no rewind]
```

The SCSI generic (sg) devices are allocated the char major number of 21. There are 256 possible SCSI generic (sg) devices:

```
[c,21,0]   /dev/sg0
[c,21,1]   /dev/sg1
....
[c,21,255] /dev/sg255
```

Note that the SCSI generic device name's use of a trailing letter (e.g. /dev/sgc) is deprecated.

Each SCSI disk (but not each partition), each SCSI CD-ROM and each SCSI tape is mapped to an sg device. SCSI devices that don't fit into these three categories (e.g. scanners) also appear as sg devices.

Pseudo devices [see the section called "Pseudo drivers"] can cause devices that are usually not considered as SCSI to appear as SCSI device names. For example an ATAPI CD-ROM may be picked up by the ide-scsi pseudo driver and mapped to /dev/scd0 .

The linux/Documentation/devices.txt file supplied within the kernel source is the definitive reference for Linux device names and their corresponding major and minor number allocations.

# Device Names in devfs

The device pseudo file system can be mounted as /dev in which case it replaces the traditional Linux device subdirectory. Alternatively it can be mounted elsewhere (e.g. /devfs) and supplement the existing device structure.

Without devfs, devices names are typically maintained in the `dev` directory of the root partition. Hence the device names (and their associated permissions) have file system persistence. The existence of a device name does not necessarily imply such a device (or even its driver) is present. To save users having to create device name entries (with the **mknod** command) most Linux distributions come with thousands of device names defined in the `/dev` directory. When applications try to open() the device name then an errno value of ENODEV indicates there is no corresponding device (or driver) currently available.

Devfs takes a different approach in which the existence of the device name is directly related to the presence of the corresponding device (and its driver).

Assuming devfs is mounted on `/dev` then SCSI devices have primary device names that might look like this:

```
/dev/scsi/host0/bus0/target1/lun0/disc    [whole disk]
/dev/scsi/host0/bus0/target1/lun0/part6   [partition 6]
/dev/scsi/host0/bus0/target1/lun0/generic [sg device for disk]

/dev/scsi/host1/bus0/target2/lun0/cd      [CD reader or writer]
/dev/scsi/host1/bus0/target2/lun0/generic [sg device for cd]

/dev/scsi/host2/bus0/target0/lun0/mt      [tape mode 0 rewind]
/dev/scsi/host2/bus0/target0/lun0/mtan    [tape mode 3 no rewind]
/dev/scsi/host2/bus0/target0/lun0/generic [sg device for tape]
```

The sg device on the third line corresponds to the "whole disk" on the first line since they have the same SCSI address (i.e. `host0/bus0/target1/lun0`). If the sg driver is a module and it has not yet been loaded (or it has been unloaded) then the "generic" device names in the above list will not be present.

[Notice the spelling of "disc" as the devfs author favours English spelling over the American variant.] It can be seen that devfs's naming scheme closely matches the SCSI addressing discussed in the section called "SCSI Addressing". It is worth noting that the IDE subsystem uses a similar devfs device naming scheme with the word "scsi" replaced with "ide". Devfs is discussed further in Chapter 12, *Devfs pseudo file system*.

# Device Names in scsidev

A utility program called **scsidev** adds device names to the `/dev/scsi` directory that reflect the SCSI address of each device. The first 2 letters of the name are the upper level SCSI driver name (i.e. either sd, sr, st or sg). The number following the "h" is the host number while the number following the "-" is meant for host identification purposes. For PCI adapters this seems to be always 0 while for ISA adapters it is their IO address. [Perhaps this field could be made more informative or dropped.] The numbers following the "c", "i" and "l" are channel (bus), target id and lun values respectively. Raw disks are shown without a trailing partition number while partitions contained within them are shown with the partition number following a "p".

The **scsidev** would typically be run as part of the boot up sequence. It may also be useful to run it after the SCSI configuration has changed (e.g. adding or removing lower level driver modules, or the use of the add/remove-single-device command). After **scsidev** has been run on my system which contains 2 disks, a cd reader and writer plus a scanner, then the following names were added in the `/dev/scsi` directory:

```
$ ls -l /dev/scsi/   # abridged
total 0
```

```
brw-------  8,   0 Sep  2 11:56 sdh0-0c0i0l0
brw-------  8,   1 Sep  2 11:56 sdh0-0c0i0l0p1
...
brw-------  8,   8 Sep  2 11:56 sdh0-0c0i0l0p8
brw-------  8,  16 Sep  2 11:56 sdh0-0c0i1l0
brw-------  8,  17 Sep  2 11:56 sdh0-0c0i1l0p1
...
brw-------  8,  24 Sep  2 11:56 sdh0-0c0i1l0p8
crw------- 21,   0 Sep  2 11:56 sgh0-0c0i0l0
crw------- 21,   1 Sep  2 11:56 sgh0-0c0i1l0
crw------- 21,   2 Sep  2 11:56 sgh1-0c0i2l0
crw------- 21,   3 Sep  2 11:56 sgh1-0c0i5l0
crw------- 21,   4 Sep  2 11:56 sgh1-0c0i6l0
br-------- 11,   0 Sep  2 11:56 srh1-0c0i2l0
br-------- 11,   1 Sep  2 11:56 srh1-0c0i6l0
```

The mapping between the SCSI generic device names (sg) and their corresponding names when controlled by other upper level drivers (i.e. sd, sr or st) can be seen by looking for name matches when the second letter is ignored. Hence "sdh0-0c0i0l0" and "sgh0-0c0i0l0" refer to the same device. By process of elimination the "sgh1-0c0i5l0" filename is the scanner since that class of devices can only be accessed via the sg interface.

The scsidev package also includes the ability to introduce names like /dev/scsi/scanner by manipulating the /etc/scsi.alias configuration file. The package also includes the useful **rescan-scsi-bus.sh** utility. For further information about **scsidev** see W6. On my system, both devfs and scsidev co-exist happily.

# Chapter 4. Kernel Configuration

The Linux kernel configuration is usually found in the kernel source in the file: `/usr/src/linux/.config`. It is not recommended to edit this file directly but to use one of these configuration options:

- **make config** - starts a character based questions and answer session

- **make menuconfig** - starts a terminal-oriented configuration tool (using ncurses)

- **make xconfig** - starts a X based configuration tool

The descriptions of these selections that is displayed by the associated help button can be found in the flat ASCII file: `/usr/src/linux/Documentation/Configure.help`

Ultimately these configuration tools edit the `.config` file. An option will either indicate some driver is built into the kernel ("=y") or will be built as a module ("=m") or is not selected. The unselected state can either be indicated by a line starting with "#" (e.g. "# CONFIG_SCSI is not set") or by the absence of the relevant line from the `.config` file.

The 3 states of the main selection option for the SCSI subsystem (which actually selects the SCSI mid level driver) follow. Only one of these should appear in an actual `.config` file:

```
CONFIG_SCSI=y
CONFIG_SCSI=m
# CONFIG_SCSI is not set
```

Some other common SCSI configuration options are:

```
CONFIG_BLK_DEV_SD          [disk (sd) driver]
CONFIG_SD_EXTRA_DEVS       [extra slots for disks added later]
CONFIG_BLK_DEV_SR          [SCSI cdrom (sr) driver]
CONFIG_BLK_DEV_SR_VENDOR   [allow vendor specific cdrom commands]
CONFIG_SR_EXTRA_DEVS       [extra slots for cdroms added later]
CONFIG_CHR_DEV_ST          [tape (st) driver]
CONFIG_CHR_DEV_OSST    [OnSteam tape (osst) driver]
CONFIG_CHR_DEV_SG          [SCSI generic (sg) driver]
CONFIG_DEBUG_QUEUES        [for debugging multiple queues]
CONFIG_SCSI_MULTI_LUN      [allow probes above lun 0]
CONFIG_SCSI_CONSTANTS      [symbolic decode of SCSI errors]
CONFIG_SCSI_LOGGING        [allow logging to be runtime selected]

CONFIG_SCSI_<ll_driver>    [numerous lower level adapter drivers]
CONFIG_SCSI_DEBUG          [lower level driver for debugging]

CONFIG_SCSI_PPA            [older parallel port zip drives]
CONFIG_SCSI_IMM            [newer parallel port zip drives]

CONFIG_BLK_DEV_IDESCSI     [ide-scsi pseudo adapter]
CONFIG_I2O_SCSI    [scsi command set over i2o bus]
CONFIG_SCSI_PCMCIA         [for SCSI HBAs on PCMCIA bus]
CONFIG_USB_STORAGE         [usb "mass storage" type]
```

```
CONFIG_MAGIC_SYSRQ          [Alt+SysRq+S for emergency sync]
                            [Alt+SyrRq+U for emergency remount ro]
```

If the root file system is on a SCSI disk then it makes sense to build into the kernel the SCSI mid level, the sd driver and the host adapter driver that the disk is connected to. In most cases it is usually safe to build the sr, st and sg drivers as modules so that they are loaded as required. If a device like a scanner is on a separate adapter then its driver may well be built as a module. In this case, that adapter driver will need to be loaded before the scanner will be recognized.

Linux distributions have many of the SCSI subsystem drivers built as modules since building all of them in would lead to a very large kernel that would exceed the capabilities of the boot loader. This leads to a "chicken and the egg" problem in which the SCSI drivers are needed to load the root file system and vice versa. The 2 phase load used by the initrd device addresses this problem (see Chapter 6, *Modules and their Parameters* for more details).

# Chapter 5. Boot Parameters

On a PC the motherboard's BIOS together with the SCSI BIOS provided by most SCSI host adapters takes care of the problem of loading the boot loader's image from a SCSI disk into memory and executing it. This may require some settings to be changed in the motherboard's BIOS. When more than one SCSI adapter is involved, the SCSI BIOS settings may need to change to indicate which one contains the disk with the boot image. The boot image make also come from an ATA (IDE) disk, a bootable CD-ROM or a floppy.

Both *lilo* and *grub* are commonly used boot loaders with Linux. Their configuration files are in `/etc/lilo.conf` and `/etc/grub.conf` [1] respectively. One difference is that after changing lilo's configuration the **lilo** command must be executed for the changes to take effect (and there is no equivalent requirement for grub). See their "man" pages for usage information. An excellent paper on lilo and the Linux bootup sequence can be found `ftp://icaftp.epfl.ch/pub/people/almesber/booting/bootinglinux-0.ps.gz` [ftp://icaftp.epfl.ch/pub/people/almesber/booting/bootinglinux-0.ps.gz]. For further information on grub see `www.gnu.org/software/grub` [http://www.gnu.org/software/grub].

Some boot parameters related to the SCSI subsystem:

```
single          [enter single user mode]
<n>             [enter run level <n> {0..6}]
root=/dev/sda6 [*]
root=/dev/scsi/host0/bus0/target0/lun0/part6 [*]
root=/dev/sd/c0b0t0u0p6    [*]
devfs=mount     [overrides CONFIG_DEVFS_MOUNT=n]
devfs=nomount  [overrides CONFIG_DEVFS_MOUNT=y]
init=<command> [executes <command> rather than init]
quiet           [reduce output to console during boot]
debug           [increase output to console during boot]
nmi_watchdog=0 [turn off NMI watchdog on a SMP machine]
max_scsi_luns=1  [limits SCSI bus scans to lun==0]
scsi_allow_ghost_devices=<n>
```

* When devfs is in use the initial read-only mount of the root partition can be done via the old /dev/sd<a><n> notation or the new devfs notation (and two of these are shown). The joint "root=/dev/sda6 single" may be useful when disk or adapter changes have broken the kernel boot load.

The "root=" argument may also be a hex number. For example, if the root partition is on `/dev/sda3` then "root=803" is appropriate. The last two digits are the minor device number discussed in an earlier section.

The default argument to the "init" parameter is `/sbin/init` (see man (8) init). If files such as `/etc/fstab` have incorrect entries, it may be useful to drop directly into a shell with "init=/bin/bash". However if shared libraries files or their paths are inappropriate this may also fail. That leaves "init=/sbin/sash" which is a statically linked shell with many useful commands (for repairing a system) built in (see man (8) sash).

When Linux fails to boot after reporting a message like:

```
  VFS: Cannot open root device 08:02
```

---

[1] One slight wrinkle with grub is that `/etc/grub.conf` is a symbolic link to `/boot/grub/grub.conf`. This can be useful to know when `/boot` is a separate partition.

then the kernel expected to find a root partition on device `/dev/sda2` and did not. The numbers in the error message are major and minor device numbers (in hex) [see the section called "Device Names" for the mapping to device names]. In such situations the "root" boot option can be useful (also the **rdev** command can be used to modify where the boot image looks for the root partition).

Lilo's configuration file `/etc/lilo.conf` can take the "root=" option in two ways. The normal way is a line like: 'root=/dev/sda2'. In this case `/dev/sda2` is converted into major and minor numbers based on the state of the system *when* the **lilo** command is executed. This can be a nuisance, especially if hardware is going to be re-arranged. The other way is a line of the form: 'append="root=/dev/sda2"' In this case the `/dev/sda2` is passed through to the kernel the next time it is started. This is the same as giving the "root=/dev/sda2" string at the kernel boot time prompt. It is interpreted by the kernel at startup (once the HBAs and their attached devices have been recognized) and thus is more flexible.

# Chapter 6. Modules and their Parameters

There are many SCSI related modules. The mid and upper level modules are listed below:

- scsi_mod.o

- sd_mod.o

- sr_mod.o

- st.o [osst.o]

- sg.o

Notice that the first 3 have "_mod" appended to their normal driver names. Lower level drivers tend to use the name (or an abbreviation) of the HBA's manufacturer (e.g. advansys) plus optionally the chip number of the major controller chip (e.g. sym53c8xx for symbios controllers based on the NCR 53c8?? family of chips).

All SCSI modules depend on the mid level. This means if the SCSI mid level is not built into the kernel and if `scsi_mod.o` has not already been loaded then a command like **modprobe st** will cause the `scsi_mod.o` module to be loaded. There could well be other dependencies, for example **modprobe sr_mod** will also cause the cdrom module to be loaded if it hasn't been already. Also if the SCSI mid level is a module, then all other SCSI subsystem drivers must be modules (this is enforced by the kernel build configuration tools).

Modules can be loaded with the **modprobe <module_name>** command which will try to load any modules that the nominated <module_name> depends on. Also <module_name> does not need the trailing ".o" extension which is assumed if not given. The **insmod <module_name>** command will also try and load <module_name> but without first loading modules it depends on. Rules for how modules can cause other modules to be loaded (with appropriate parameters appended) are usually placed in the file `/etc/modules.conf`. [Note that in earlier Linux kernels this file was often called `/etc/conf.modules`.] For further information about the format of this file try **man modules.conf**.

Any module can have its allowable command line parameters queried with this command: **modinfo -p <module_name>**.

When upper level drivers are initialized and if there are no hosts active then the mid level will attempt to load a module called "scsi_hostadapter". An "alias" can then be used to associate "scsi_hostadapter" with the actual name of the lower level (adapter) driver. For example, a line like "alias scsi_hostadapter aic7xxx" in the `/etc/modules.conf` file would cause the aic7xxx module to be loaded (if there were no lower level drivers already active). [1]

There is a special relationship between the module parameter "scsi_hostadapter" and the initrd file system. For more information see **man initrd** and **man mkinitrd**. [2]

---

[1] There is a sequencing issue here if the root file system is on the SCSI device controlled by the lower level (adapter) driver to be loaded since it contains the `/etc/modules.conf` file. Also there is the issue of how the boot loader obtains the initial kernel image from a SCSI device (often from the (master) boot record). The latter is usually taken care of by the system's or adapter card's BIOS.

[2] An example of using **mkinitrd**: assume the root partition is on a SCSI disk connected to a controller from Adaptec that requires the aic7xxx driver. After building a kernel with the aic7xxx driver specified as a module then load the image into the normal place (probably in the `/boot` directory). Next make sure a line like "alias scsi_hostadapter aic7xxx" is in the `/etc/modules.conf` file. Then from the `/boot` directory execute a

---

line like **mkinitrd /boot/initrd-2.4.5.img 2.4.5** (this assumes lk 2.4.5 is being build). This should result in the file `initrd-2.4.5.img` being created. The `/etc/lilo.conf` should then have a section added looking something like this:

```
image=/boot/vmlinuz-2.4.5
    label=linux
    initrd=/boot/initrd-2.4.5.img
    read-only
    root=/dev/sda7
```

The following should also be selected in the kernel configuration:

```
CONFIG_BLK_DEV_RAM=y
CONFIG_BLK_DEV_RAM_SIZE=4096
CONFIG_BLK_DEV_INITRD=y
```

See also `Documentation/initrd.txt`.

# Chapter 7. Proc pseudo file system

The proc pseudo file system provides some useful information about the SCSI subsystem. The kernel configuration option that selects "proc_fs" is CONFIG_PROC_FS and in almost all cases it should be selected. SCSI specific information is found under the directory `/proc/scsi`. Probably the most commonly accessed entry is **cat /proc/scsi/scsi** which lists the attached SCSI devices. See the section called "proc interface" for more details.

The lower level drivers are allocated proc_fs entries of the form:

```
/proc/scsi/<driver_name>/<scsi_adapter_number>
```

where the <driver_name> is something like "aic7xxx" or "BusLogic". The <scsi_adapter_number> (also known as the host number) is the same number that was discussed in the section called "SCSI Addressing". Note that one driver may control one or more hosts. What is stored in this file is lower level driver dependent (and in the case of some adapter drivers it is possible to set parameters via this file). When reporting problems to newsgroups or maintainers it is useful to include the output of this file (e.g. **cat /proc/scsi/aic7xxx/0** ).

The cdrom driver provides information about attached cdrom devices in the `/proc/sys/dev/cdrom` directory. This will include both SCSI devices (i.e. those controlled by the sr driver) and IDE devices (i.e. those controlled by the ide-cd driver). See the section called "sr proc interface".

The sg driver provides information about its state and attached hosts and devices in the `/proc/scsi/sg` directory. See the section called "sg proc interface".

More general information on the proc pseudo file system can be found in the kernel source file: `/usr/src/linux/Documentation/filesystems/proc.txt`.

# Chapter 8. Mid Level, Unifying layer

The SCSI mid level is common to all usage of the SCSI subsystem. Probably its most important role is to define internal interfaces and services that are used by all other SCSI drivers. These internal mechanisms are not discussed in this document [see ref: W2].

The primary kernel configuration parameter "CONFIG_SCSI" determines whether the mid level is built in (when "=y") or a module (when "=m"). If "CONFIG_SCSI=m" then all other SCSI subsystem drivers must also be modules.

When the mid level is built as a module then it probably never needs to be loaded explicitly because using 'modprobe' to load any other SCSI subsystem module will cause the mid level to be loaded first (if it is not already).

Some upper and lower level drivers do not (fully) load if there are no devices for that driver to control. Sometimes the report is loud as in this case for the imm driver which controls zip drives connected to a parallel port:

```
$ modprobe imm
    imm.o: init_module: No such device
```

**lsmod** will not show the "imm" module as loaded. In other cases the result is more subtle. For example, if the sg driver is loaded in a system with no (real or pseudo) scsi devices then the `/proc/scsi/sg` directory will not appear. [It will be created when the first scsi device is recognized.]

# boot parameters

SCSI drivers that are built into the kernel are checked in a pre-determined order to see if HBAs that they can control are present. The user has no control over this order which in most cases is arbitrary but in the case of some older ISA adapters is required to stop misidentification [1] .

```
scsi_logging=<n>
    where <n> is 0 to turn logging off
    where <n> is non-zero to turn logging on

max_scsi_luns=<n>
    where <n> is a number between 1 and 8 (< lk 2.4.7),
     >= lk 2.4.7 the upper limit can be much larger

scsi_allow_ghost_devices=<n>
    where (<n> - 1) is the maximum lu number to ghost if the
        the corresponding device is offline. When <n>==0
        (default) then don't ghost any devices (in lk 2.4.26
        and later)

scsihosts=host0:hosts1::host3
```

---

[1] PCI adapters are much "safer" for initialization code than the older ISA adapters. Hence the order of initialization of PCI adapters is unlikely to lead to lockups. In this case the order of initialization (and thus SCSI adapter numbers) of built in drivers may be modified by changing the order of entries in the SCSI subsystem Makefile ( `/usr/src/linux/drivers/scsi/Makefile`). Beware: some adapters may be recognized by more than one lower level driver (e.g. those based on NCR chipsets).

The recently introduced devfs defines a "scsihosts" boot time parameter to give the user some control over this. See the devfs documentation [ref: W5] for a description. The host names given in the list to the "scsihosts" boot option are the names of lower level drivers (e.g. "scsihosts=advansys:imm::ide-scsi"). [2] [3] Devfs does not need to be present for "scsihosts" to be used. The "scsihosts" parameter, if given, is echoed during in the boot up messages. For example:

```
scsi: host order: advansys:imm::ide-scsi
```

Also if multiple HBA are present in a system then they are scanned in a fixed order (see footnote). The "scsihosts" parameter only effects how these HBAs are indexed (i.e. which SCSI adapter numbers are associated with them by the kernel). In the above example, if the "imm" driver is not found during boot up, then the scsi adapter number "1" is not allocated. If the "imm" driver is later loaded as a module, then it will adopt scsi adapter number "1". If a driver that is not named in "scsihosts" is found, then it will get the next available scsi adapter number (e.g. a built in aic7xxx driver would get scsi adapter number "2" in the above example).

A full list of kernel parameters with some explanations can be found in the file `/usr/src/linux/Doc-umentation/kernel-parameters.txt` .

# module parameters

If SCSI disks are present in the system then it usually is better to build the mid level driver into the kernel. However if the SCSI subsystem is only being used periodically (e.g. to burn CD-Rs on an ATAPI CD writer) then building the mid level as a module is fine. The module load time options are the same as the driver's built in options:

```
scsi_logging_level=<n>
    where <n> is the logging level mask (0 for logging off)
max_scsi_luns=<n>
scsihosts=host0::host2
scsi_allow_ghost_devices=<n>
```

# proc interface

To display the SCSI devices currently attached (and recognized) by the SCSI subsystem use **cat /proc/scsi/scsi**.

The output looks like this:

```
    Attached devices:
    Host: scsi0 Channel: 00 Id: 02 Lun: 00
      Vendor: PIONEER  Model: DVD-ROM DVD-303  Rev: 1.10
      Type:   CD-ROM               ANSI SCSI revision: 02
    Host: scsi1 Channel: 00 Id: 00 Lun: 00
```

---

[2] Either comma or colon can be delimiters for "scsihosts". This means that "scsihosts=advansys,imm,,ide-scsi" is also valid. Also if a machine's boot sequence involves an "initrd" stage (look in `/etc/grub.conf` or `/etc/lilo.conf` to find out if this is the case), then the **mkinitrd** command should be run after a change to the "scsihosts" boot time parameter. This will generate a new initrd image that needs to be put in the correct place (most probably in the `/boot` directory).

[3] Using "scsihosts" can lead to a situation in which the computer's BIOS finds the boot track (and hence boot time parameters set in lilo or grub) on one disk while the kernel finds the root partition on another disk. This can be quite confusing when it is unplanned. Hence after changing (or adding) "scsihosts" in lilo or grub's configuration, it may be wise to boot the machine to see which disks are accessed.

```
        Vendor: IBM       Model: DNES-309170W     Rev: SA30
        Type:   Direct-Access       ANSI SCSI revision: 03
```

After the "Attached devices:" line there are 3 lines for each recognized device. The first of these lines is SCSI address information discussed in the section called "SCSI Addressing". The following 2 lines of data are obtained from a INQUIRY command that was performed on the device when it was attached. See the section called "Generic driver (sg)" for the relationship between the ordering of these devices compared with the sg driver's ordering (which most of the time is the same).

Existing devices can be removed using **echo "scsi remove-single-device <h> <b> <t> <l>" > /proc/scsi/scsi** where the variables are host, bus (channel), target (scsi id) and lun. The success (or otherwise) of this command can be determined by sending a subsequent **cat /proc/scsi/scsi** command. The removal will fail if the device is busy (e.g. if a file system on the device is mounted).

New devices can be added using **echo "scsi add-single-device <h> <b> <t> <l>" > /proc/scsi/scsi** where the variables are host, bus (channel), target (scsi id) and lun. The success (or otherwise) of this command can be determined by sending a subsequent **cat /proc/scsi/scsi** command. [4]

The SCSI subsystem does not support hot-plugging of SCSI devices (there may also be electrical issues on the associated SCSI parallel bus). It is recommended that those who use add+remove-single-device make sure that other devices on that SCSI bus are inactive if re-plugging is going to take place.

To output a list of internal SCSI command blocks use **echo "scsi dump <n>" > /proc/scsi/scsi** where the numeric value of <n> doesn't matter. This is probably only of interest to people chasing down bugs within the SCSI subsystem.

To start (or stop) logging information being sent to the console/log use **echo "scsi log <token> <n>" > /proc/scsi/scsi** where <token> is one of: {all, none, error, timeout, scan, mlqueue, mlcomplete, llqueue, llcomplete, hlqueue, hlcomplete, ioctl} and <n> is a number between 0 and 7. The tokens "all" and "none" don't take an <n> argument. Prefix meanings:

```
   hl    upper level drivers [exception: sg uses "timeout"]
   ml    mid level
   ll    lower level drivers
          [adapter drivers often have there own flags]
```

The value "0" turns off logging while "7" maximizes the volume of output. Logging information will only be output if CONFIG_SCSI_LOGGING was selected in the kernel build.

## Warning

Warning: "scsi log all" (and several other variants) can cause a logging infinite loop if the log file (typically /var/log/messages ) lies on a SCSI disk. Either turn off the kernel logging daemon or direct its output to a non SCSI device.

---

[4] The parsing of "add-single-device" and "remove-single-device" is rather inflexible. Hence it is best to stay close to the demonstrated syntax with no extra spaces (and no tabs).

# Chapter 9. Upper level drivers

The upper level drivers maintain the kernel side of the OS interface for the logical class of devices they represent (e.g. disks). They are also responsible for managing certain kernel and SCSI subsystem resources such as kernel memory and SCSI command structures. Applications in the user space access these drivers by opening a special file (block or char) typically found in the `/dev` directory tree.

## Disk driver (sd)

Two types of SCSI devices are accessible via the sd driver:

- "direct access" devices which are usually magnetic disks. [SCSI peripheral device code is 0]

- "Optical memory devices" which are often called MOD disks. [SCSI peripheral device code is 7]

The sd driver is a block device which means that it is closely associated with the block subsystem. It also supports the concept of partitions. [**man sd** dates from 1992.]

The sd driver is capable of recognizing 128 disks when it is loaded at kernel boot time or later as a module. However, once it is loaded, it will only recognize a fixed number of additional disks. The number of additional disks that can be accommodated is set by the kernel configuration parameter CONFIG_SD_EXTRA_DEVS whose default value is 40.

## sd boot parameters

None.

## sd module parameters

The sd driver takes no parameters when loaded as a module. Note that its module name is `sd_mod.o`.

## CDROM driver (sr or scd)

CDROM and DVD drives (and WORM devices) are accessible via the sr upper level device driver. While "sr" is the device driver name, "sr_mod" is its module name. The device file name is either `/dev/sr<n>` or `/dev/scd<n>`.

Following is a diagram illustrating the CDROM subsystem of which sr is a part:

The architecture of the CD-ROM subsystem.

This diagram glosses over some of the differences between the protocol stacks. CDROM device names are *not* maintained by the uniform CDROM layer but rather by each individual protocol stack. In the case of the SCSI subsystem, device names are maintained by the sr driver while the IDE subsystem maintains device names with its central "ide" driver (i.e. not by the ide-cd driver). USB and IEEE1394 cd devices names are maintained by their respective stacks. This may partially explain why the `/dev/cdrom` is often a symbolic link to the appropriate subsystem's device name.

Two types of SCSI devices are accessible via the sr driver:

- CD-ROM devices (including DVD players) [SCSI peripheral device code is 5]

- "Write-once read-multiple" devices which are known as WORMs. [SCSI peripheral device code is 4]

The sr driver is capable of recognizing 256 CDROM/DVD drives when it is loaded at kernel boot time or later as a module. However, once it is loaded, it will only recognize a fixed number of additional drives. The number of additional drives that can be accommodated is set by the kernel configuration parameter CONFIG_SR_EXTRA_DEVS whose default value is 2.

People often use the **dd** command to read a data CDROM containing an iso9660 file system. If a count argument is not given then the **dd** command will read the number of 2048 byte sectors indicated by the SCSI Read Capacity command. Unfortunately this can include unwritten (or "run out") sectors at the end of the image that can cause I/O errors. Use the **isosize** command (see its man page) to find the true length of the iso9660 image and use that in the "count=" argument given to the **dd** command.

# sr boot parameters

None.

# sr module parameters

Doing a test to find out if a cdrom drive supports XA mode (mode 2) triggers firmware bugs on some drives. Consequently the check for XA mode support is turned off by default. The following module parameter is provided:

```
xa_test=<0|1>
```

to override the default. [Currently there seems to be no way to turn on XA mode testing when the sr driver is built into the kernel.]

# sr proc interface

All the following files are readable by all and produce ASCII output when read:

```
/proc/sys/dev/cdrom/autoclose
/proc/sys/dev/cdrom/autoeject
/proc/sys/dev/cdrom/check_media
/proc/sys/dev/cdrom/debug
/proc/sys/dev/cdrom/info
/proc/sys/dev/cdrom/lock
```

They reflect the current state of the CDROM subsystem. This location is part of the procfs's window through to the sysctl configuration mechanism (see **man sysctl**). All but info are writable by the superuser. There is a column for each CDROM and DVD player in the system in info (not just SCSI devices).

As an example, the auto eject feature can be turned on by the superuser with the command **echo "1" > /proc/sys/dev/cdrom/autoeject**. This will cause cdroms to be ejected from the drive when unmounted.

# ATAPI cdroms

Many Linux users have no SCSI devices (or adapters) in their systems. They become a little perplexed as to why cd writer software (e.g. **cdrecord** and **cdrdao**) and cd music reading programs (e.g. **cdparanoia**) use the Linux SCSI subsystem. The answer is that these programs need lower level access to these devices.

ATAPI (ATA Packet Interface) is essentially a SCSI command set sent over an ATA [1] transport. [The discussion in this section is also applicable to ATAPI tape drives and ATAPI floppy drives.]

Currently both **cdrecord** and **cdparanoia** interface to the SCSI generic driver (sg) and, in the case of ATAPI cd devices, use the ide-scsi pseudo device driver to access the hardware. This may change in the future as in the 2.4 series kernels a packet interface ioctl has been added to the uniform cdrom layer (see the diagram in the section called "CDROM driver (sr or scd)" above). [2]

The default action of the IDE subsystem in Linux is to claim all ATA devices for its built-in drivers. In the case of an ATAPI cd writer, it will be claimed by the built-in ide-cd driver. Once this has happened, the SCSI subsystem is unable to get control over an ATAPI device. The ide-scsi (pseudo lower level SCSI) driver can only register ATAPI devices in the SCSI subsystem that have *not* already been claimed by the IDE subsystem.

Notice the *built-in* qualification in the previous paragraph. If both the ide-cd and ide-scsi drivers are modules then the first one loaded will claim the ATAPI cd devices (e.g. cd/dvd readers and writers). Furthermore you can switch the controlling driver module by **rmmod**-ing one and **modprobe**-ing the other.

Probably the most flexible way to instruct the IDE core driver that you want the cd writer at `/dev/hdd` accessible to **cdrecord** is to use the kernel boot option: "hdd=ide-scsi". This will cause the ide-cd driver to bypass `/dev/hdd` (irrespective of whether ide-cd driver is built-in or a module). As long as the ide-scsi driver is built-in or a module then it will "capture" the cd writer at `/dev/hdd` (with the IDE core driver loading the ide-scsi module if required).

The ide-cd driver module can be instructed to ignore certain ATA devices with the following syntax:

```
modprobe ide-cd ignore='hdc hdd'
```

In this case the ide-cd driver will ignore the devices at `/dev/hdc` and `/dev/hdd`. This effect can also be accomplished by placing a line like this: "options ide-cd ignore=hdd" in the `/etc/modules.conf` file.

A new option added in the lk 2.4 series is of the form "hdd=scsi". This option seems to have a similar function to the "hdd=ide-scsi" option discussed above. Furthermore "hdd=scsi" can only be used if both the SCSI mid-level and the ide-scsi drivers are built into the kernel (otherwise "BAD OPTION" is reported by the ide_setup function).

To find out whether an ATAPI cd device is "owned" by the SCSI subsystem, the output of **cat /proc/scsi/scsi** can be checked. Another technique is to observe the "drive name:" line of **cat /proc/sys/dev/cdrom/info** for "sr" entries. The following output is from my system:

```
$ cat /proc/sys/dev/cdrom/info
CD-ROM information, Id: cdrom.c 3.12 2000/10/18

drive name:            sr1     sr0
drive speed:           16      0
drive # of slots:      1       1
Can close tray:        1       1
Can open tray:         1       1
Can lock tray:         1       1
```

---

[1] ATA is the modern name for what was previously known as IDE and/or EIDE. Note that the subsystem that controls ATA devices in Linux is called the "IDE" subsystem for historical reasons.

[2] Other ATA devices such as tapes and floppies often use the ATAPI interface. However, the vast majority of ATA disks do *not* use the ATAPI interface.

```
Can change speed:        1        1
Can select disk:         0        0
Can read multisession:   1        1
Can read MCN:            1        1
Reports media changed:   1        1
Can play audio:          1        1
Can write CD-R:          1        0
Can write CD-RW:         1        0
Can read DVD:            0        1
Can write DVD-R:         0        0
Can write DVD-RAM:       0        0
```

Once an ATAPI cd writer at /dev/hdd has been registered by the SCSI subsystem, then cdroms should be mounted via the "scd" device name and cd players should also use the "scd" device. Strangely the **hdparm** command should still use the   /dev/hdd device file (or the "echo ... > /proc/ide/hdd/settings" method described in this section). [3]

# Tape driver (st)

The tape driver interface is documented in the file /usr/src/linux/drivers/scsi/README.st and on the st(4) man page (type **man st**). The file README.st also documents the different parameters and options of the driver together with the basic mechanisms used in the driver.

The tape driver is usually accessed via the **mt** command (see **man mt**). **mtx** is an associated program for controlling tape autoloaders (see mtx.sourceforge.net [http://mtx.sourceforge.net]).

The st driver detects those SCSI devices whose peripheral device type is "Sequential-access" (code number 1) unless they appear on the driver's "reject_list". [Currently the OnStream tape drives (described in a following section) are the only entry in this reject_list.]

The st driver is capable of recognizing 32 tape drives. There are 8 device file names for each tape drive: a rewind and non-rewind variant for each of 4 modes (numbered 0 to 3). See the tape device file name examples in the section called "Device Names" on device names. Any number of tape drives (up to the overall limit of 32) can be added after the st driver is loaded.

ATAPI tape drives can be controlled by this driver with help from the ide-scsi pseudo adapter driver. The discussion in the section called "ATAPI cdroms" also applies for ATAPI tape drives (and ATAPI floppies).

## st boot parameters

```
st=xxx[,yyy] where xxx is one of the following:
```

---

[3] In the linux 2.4 kernel series there has been an increase in problems when the ide-scsi driver is used so that **cdrecord** can control ATAPI (IDE) cd writers. The problem may be related to the aggressive manner in which the IDE subsystem attempts to optimize the speed of data transfers on devices it controls. Some people experiencing timeouts and machine lockups have found that reducing the DMA setting via the **hdparm** command has fixed the problem. If the cd writer is connected to /dev/hdd then users have reported success with one of these two commands:

```
hdparm -d0 -c1 /dev/hdd
hdparm -d 1 -X 34 /dev/hdd
```

The first one turns off DMA completely while the second one sets it in "multiword DMA mode 2". Cd writers do not need the types of speeds that modern disks utilize. Even burning at "x16" implies a sustained transfer rate of 16 times 150 KB/sec which is approximately 2.4 MB/sec, not really that fast. There has also been a report that moving a cd writer off a high speed IDE controller (Promise) and back to the motherboard's lower speed IDE controllers has fixed a random IDE bus reset problem. Another report suggests reducing (or turning off) the DMA on the IDE hard disk can also stop lockups.

```
buffer_kbs:<n>
write_threshold_kbs:<n>
max_buffers:<n>
max_sg_segs:<n>

(The old boot parameters st=aa[,bb[,cc[,dd]]] supported but deprecated)
```

The default driver buffer size (buffer_kbs) is 32 (i.e. 32 KB). The default asynchronous write threshold (write_threshold_kbs) is 30 (i.e. 30 KB). The default number of buffers allocated at initialization (max_buffers) is 4. The default number of scatter/gather segments to use (max_sg_segs) is 32.

## st module parameters

```
buffer_kbs=<n>
write_threshold_kbs=<n>
max_buffers=<n>
max_sg_segs=<n>
```

## st proc interface

None.

## osst driver for OnStream devices

There is an auxiliary tape driver for tape drives manufactured by OnStream. It is an additional upper level driver and can co-exist with the st driver. Its driver name is "osst" (as is its module name).

The OnStream SC-x0 SCSI tape drives can not be driven by the standard st driver, but instead need this special osst driver and use the /dev/osst<x> char device nodes (major 206). [Where <x> follows the same naming scheme as st devices outlined in the section called "Device Names".] Via usb-storage and ide-scsi, you may be able to drive the USB-x0 and DI-x0 drives as well. Note that there is also a second generation of OnStream tape drives (ADR-x0) that supports the standard SCSI-2 commands for tapes (QIC-157) and can be driven by the standard driver st. For more information, you may have a look at the kernel source file /usr/src/linux/drivers/scsi/README.osst. More info on the OnStream driver may be found on linux1.onstream.nl/test/ [http://linux1.onstream.nl/test/].

# Generic driver (sg)

All types of SCSI devices are accessible via the sg driver. This means devices such as CDROM drives can be accessed both via the sr and sg drivers. Other SCSI devices such as scanners can only be accessed via the sg driver. The sg driver is capable of recognizing 256 SCSI devices. Any number of devices (up to the overall limit of 256) can be added after the sg driver is loaded.

See reference W4 for the SCSI Generic (sg) driver documentation (also found there is the sg_utils package). For SCSI standards see reference W1 and for a book on the subject of SCSI programming and pass through mechanisms see reference B3.

The sg driver in lk 2.4 is "version 3" which adds an additional interface structure and some new ioctl()s. The most interesting new ioctl() is SG_IO which sends a SCSI command and waits for its response. See the Linux Documentation Project site: www.tldp.org/HOWTO/SCSI-Generic-HOWTO/ [http://www.tldp.org/HOWTO/SCSI-Generic-HOWTO/] for a full description of the sg driver. A (possibly lat-

er) version of this document can be found at `www.torque.net/sg/p/sg_v3_ho.html` [http://www.torque.net/sg/p/sg_v3_ho.html].

The abbreviation "sg" is used within the kernel to refer both to the SCSI generic driver and the scatter-gather capability offered by many modern IO devices (usually associated with DMA). The context usually makes it clear which one is being referred to. As an example, note the contorted sg ioctl() named SG_GET_SG_TABLESIZE where the second "SG" refers to scatter gather.

The public interface for sg is found in the file: `/usr/src/linux/include/scsi/sg.h`. Depending on the distribution this may or may not contain the same information as `/usr/include/scsi/sg.h` which is controlled by the GNU library maintainers. If these 2 files are not the same use the former header file. Those writing applications based on sg should see its documentation for more on this matter.

The sg driver registers all SCSI devices (with a current maximum of 256) as they are seen. Each newly registered SCSI device gets allocated the next available minor device number. At least initially this will be the same sequence that devices are displayed in mid level's **cat /proc/scsi/scsi**. The sg devices device mapping can be seen with **cat /proc/scsi/sg/devices** or **cat /proc/scsi/sg/device_strs**. Differences between **cat /proc/scsi/scsi** and sg orderings will appear when a low level driver is removed (e.g. **rmmod aha1542**) or when a device is removed with remove-single-device. The sg driver will leave remaining SCSI device mapping to minor device numbers unchanged. This potentially leaves a "hole" in the sg mapping. An example follows:

```
$ cat /proc/scsi/scsi
Attached devices:
Host: scsi0 Channel: 00 Id: 00 Lun: 00
  Vendor: IBM      Model: DNES-309170W     Rev: SA30
  Type:   Direct-Access            ANSI SCSI revision: 03
Host: scsi1 Channel: 00 Id: 02 Lun: 00
  Vendor: PIONEER  Model: DVD-ROM DVD-303  Rev: 1.10
  Type:   CD-ROM                   ANSI SCSI revision: 02
Host: scsi1 Channel: 00 Id: 06 Lun: 00
  Vendor: YAMAHA   Model: CRW4416S         Rev: 1.0g
  Type:   CD-ROM                   ANSI SCSI revision: 02

$ cat /proc/scsi/sg/device_strs
IBM             DNES-309170W             SA30
PIONEER         DVD-ROM DVD-303          1.10
YAMAHA          CRW4416S                 1.0g

$ echo "scsi remove-single-device 1 0 2 0" > /proc/scsi/scsi

$ cat /proc/scsi/scsi
Attached devices:
Host: scsi0 Channel: 00 Id: 00 Lun: 00
  Vendor: IBM      Model: DNES-309170W     Rev: SA30
  Type:   Direct-Access            ANSI SCSI revision: 03
Host: scsi1 Channel: 00 Id: 06 Lun: 00
  Vendor: YAMAHA   Model: CRW4416S         Rev: 1.0g
  Type:   CD-ROM                   ANSI SCSI revision: 02

$ cat /proc/scsi/sg/device_strs
IBM             DNES-309170W             SA30
<no active device>
```

```
YAMAHA          CRW4416S              1.0g
```

Notice how the sg driver maintains the row positions of the remaining devices in the "device_strs" output. So when the Pioneer dvd player is removed, a hole opens up in the sg device mapping which is not reflected in the **cat /proc/scsi/scsi** output. That "hole" corresponds to the device name `/dev/sg1`.

The new sg_io_hdr interface includes a data transfer residual count field called "resid". Only some lower level adapters support this feature and those that don't always yield zero in this field. At the time of writing the advansys, aha152x and the sym53c8xx drivers support this feature.

# sg boot parameters

The sg driver maintains a reserved buffer for each open file descriptor. The purpose is to guarantee applications that data transfers up to the size of the reserved buffer will not fail for lack of kernel memory. This is important for applications like cdrecord that cannot easily recover (the CDR) from a ENOMEM error.

In the absence of the boot parameter 'sg_def_reserved_size' or the sg module parameter 'def_reserved_size', then each time a sg file descriptor is opened the reserved buffer size is inherited from SG_DEF_RESERVED_SIZE which is defined in `include/linux/sg.h`.

The SG_DEF_RESERVED_SIZE define value can be overridden by this kernel boot option:

```
sg_def_reserved_size=<n>
```

# sg module parameters

When the sg module is loaded the SG_DEF_RESERVED_SIZE define value can be overridden by supplying this option:

```
def_reserved_size=<n>
```

# sg proc interface

All the following files are readable by all and produce ASCII output when read. The file 'def_reserved_size' is also writable by root. The ASCII output has been formatted in such a way as to be human and machine readable (and hence a compromise). Use Unix commands of the form **cat device_hdrs devices** to see the output of tables.

```
/proc/scsi/sg/debug       [internal state of sg driver]
/proc/scsi/sg/def_reserved_size
                          [like boot/module load parameter]
/proc/scsi/sg/devices     [table of numeric device data]
/proc/scsi/sg/device_hdr  [column headers for sg/devices]
/proc/scsi/sg/device_strs [table of strings from INQUIRY]
/proc/scsi/sg/hosts       [table of numeric host data]
/proc/scsi/sg/host_hdr    [column headers for sg/hosts]
/proc/scsi/sg/host_strs   [table of string ids for hosts]
/proc/scsi/sg/version     [sg version number and date]
```

All the above files are owned by root and readable by all while `def_reserved_size` is writable by root. For the `devices` and `device_strs` files the first row output corresponds to /dev/sg0 (sg

minor device number 0). The second row output corresponds to `/dev/sg1`, etc. For the `hosts` and `host_strs` files the first row output corresponds to host (adapter number) 0, etc. For numeric tables a missing device or host is indicated by a row of "-1" values. For string tables a missing device or host is indicated by a row containing "<no active device/host>".

# Chapter 10. Lower Level drivers

There are too many SCSI low level drivers to detail in this document. As an alternative to giving any superficial overview here, the reader is given suggestions of places to look for further information.

The source directory for the SCSI subsystem in the Linux kernel is a good place to start: `/usr/src/lin-ux/drivers/scsi`. Several drivers have information in a "readme" file: `README.<driver_name>`. Others have extensive information at the top of their ".c" file This information often includes a version number, change logs and kernel boot time and module load time options. Often the latter information can be found in the installation guides of the various Linux distributions. Sometimes the driver maintainer will have a web site containing the most recent bug fix information. Official maintainers are listed in the `/usr/src/linux/MAINTAINERS` file. If there is nothing there, look in the relevant ".c" file in the SCSI subsystem directory. Some old drivers have no active maintainers. In such cases posting to the linux-scsi newsgroup may help [see N1 ].

For an overview of the drivers supplied with the kernel source tree, use one of the kernel configuration programs (e.g. **cd /usr/src/linux; make menuconfig**). The help information associated with each selection can be found together in one (large) flat file at `/usr/src/linux/Documenta-tion/Configure.help`. Drivers can be obtained from other places. It is unlikely that a SCSI driver made for the lk 2.2 series (or before) will build or operate successfully in the lk 2.4 series. [From a programmatic viewpoint there are not a lot of things that need changing.] Drivers may even be only available in binary form, in which case make sure that you trust the provider and follow their instructions closely.

Lower level drivers can support either of 2 error handling strategies. The older one is considered obsolete while the newer one is often called "new_eh". The advantage of "new_eh" is that it uses a separate kernel thread per host (named "scsi_eh_<n>" where <n> is the host number) to facilitate error recovery. Both error handling strategies were also available in the lk 2.2 series in which very few adapter drivers used "new_eh". In the lk 2.4 series, more drivers are using it and the plan for the forthcoming lk 2.5 development series is to drop mid level support for the older, obsolete error strategy.

Drew Eckhardt's SCSI-HOWTO document [see reference W7 ] goes into much more detail about lower level (adapter) drivers than this document. Since that SCSI-HOWTO is 5 years old, many things have changed and more drivers have been added.

There is a lower level driver called *scsi_debug* that simulates one or more "direct access" devices (i.e. disk(s)) using the computer's memory. From lk 2.4.17 it acts as a "ram disk". While there are many ram disk implementations available in Linux (e.g. ramfs), scsi_debug may help to isolate a defective scsi driver in a problematic installation. See `scsi_debug.c` for further information.

# Pseudo drivers

SCSI can be viewed as a command set and a set of hardware buses that convey that command set. Those hardware buses can be further divided into those used exclusively for SCSI (e.g. ultra wide), those shared with other protocols (e.g. USB, IEEE 1394) and those buses not defined by the various SCSI standards. In the final category there are several interesting examples including ATAPI CD writers and PC parallel bus ZIP drives. Such devices use the SCSI command set (or something very close to it) over a foreign bus.

This section briefly outlines various pseudo lower level drivers which essentially communicate with other Linux subsystems in order to send the SCSI command set to devices controlled by those other subsystems. This raises some ownership issues that often confuse users and result in many questions to the maintainers.

**IDE-SCSI.** From configuration point of view, ide-scsi will grab and try to control every ATA (a.k.a. IDE) device which doesn't have a "native" driver attached (such as ide-cd, ide-tape, etc). So for example,

if both ide-cd and ide-scsi are compiled into the kernel in a system which has an ATAPI cdrom, ide-cd will get to control it. If only ide-scsi is compiled in, it will get the device. There are some kernel boot time parameters to control which driver gets which device.

The preferences of the IDE subsystem can be overridden with one of these kernel boot time parameters (of which the first is most interesting for this subsystem):

- hdx=ide-scsi

- hdx=ide-cdrom

- hdx=ide-floppy

[The term *hdx* is used to refer to one of the IDE/ATA devices in {hda, hdb, hdc ...}.] In the 2.4 series "hdx=scsi" was added but it is not very useful, see see the section called "ATAPI cdroms".

When the driver is running, the device will be accessible using the SCSI device (`/dev/sda`, `/dev/sr0` , etc), and not through the corresponding `/dev/hdx` device. Still, the `/dev/hdx` device will be available, but only for configuration.

All the generic IDE configuration parameters (DMA on/off, 32-bit I/O, unmasking irq's, etc) are available by using the `/dev/hdx` device, for example to enable DMA:

```
hdparm -d1 /dev/hdx
```

[1] Using **cat /proc/ide/hdx/settings** will show the available settings. All the generic IDE driver settings will be available there, as well as the following "ide-scsi specific" settings:

- bios_cyl

- bios_head

- bios_sect

- transform

- log

The first three choose the virtual geometry that the drive will return to the sd driver, in case it's a disk drive (ZIP, etc). "transform" will configure/enable/disable the SCSI to ATAPI CDB transformation layer:

- bit 0: Enable(1)/Disable(0) transformation for commands not originated from the sg driver.

- bit 1: Enable/Disable transformation for commands issued using the sg driver.

"log" will log debugging information. This is useful also to debug user-space programs using the sg driver, as it will list the CDB traffic on the bus -- each issued command, along with its completion status. To enable/disable a specific settings, use something like:

```
echo "log:1" > /proc/ide/hdx/settings
```

To turn off the "using_dma" flag use:

---

[1] It has been reported that in some distributions the attempt to use the hdparm command fails. In this case use the "echo ... > /proc/ide/hdx/settings" form.

```
echo "using_dma:0" > /proc/ide/hdx/settings
```

**PPA + IMM.**   Iomega ZIP drives come in a variety of flavours including parallel port, SCSI, and ATAPI. The parallel port versions (both old and new) are driven by ppa and imm respectively.

The parallel port ZIP drives are actually SCSI devices which tunnel SCSI commands over the parallel port using interfaces called VPI0 (older-style) and VPI2 (newer-style). The ppa driver is the VPI0 host implementation and the imm driver is the VPI2 host implementation.

The way it works is that the HBA is a chip inside the ZIP drive, so that the host adapter and the peripheral are in the same actual case.

**PPSCSI.**   The new, not-yet-integrated, architecture for devices that use SCSI over a parallel port cable is ppscsi. The ppscsi module provides the boiler plate code and makes it easy to write implementations for different interfaces.

Each ppscsi protocol module registers itself with the ppscsi module, passing in a list of entry points for the various things that are common to all protocol drivers.

The structure of the PPSCSI drivers.

The plan is that the ppscsi architecture will absorb both the ppa and imm drivers and protocol modules; only vpi0 has been written so far. See `www.torque.net/parport/ppscsi.html` [http://www.torque.net/parport/ppscsi.html].

**USB.**   USB classifies a group of devices as "mass storage" (e.g. disks) and interacts with these using the SCSI command set. The module name is "usb-storage". See `www.one-eyed-alien.net/~md-harm/linux-usb` [http://www.one-eyed-alien.net/~mdharm/linux-usb].

There is also the usb/microtek driver for controlling X6 USB scanners from Microtek. When configured, the SANE application uses the sg driver to send SCSI commands over USB to control this scanner.

**I2O.**   See kernel source file `/usr/src/linux/drivers/i2o/io2_scsi.c` .

**IEEE 1394.**   Support for IEEE 1394 devices that use the SBP-2 protocol is now available (lk 2.4.7). See the IEEE 1394 paragraph in this section for more information.

**iSCSI.**   An IETF draft is taking shape for iSCSI. This sends the SCSI command set over a TCP network connection. iSCSI seems to be gaining popularity quickly and there are several implementations for Linux taking shape. One implementation is at `sourceforge.net/projects/intel-iscsi/` [http://sourceforge.net/projects/intel-iscsi/]. Use your favourite search engine to find other projects.

# Chapter 11. Raw devices

A raw device can be bound to an existing block device (e.g. a disk) and be used to perform "raw" IO with that existing block device. Such "raw" IO bypasses the caching that is normally associated with block devices. Hence a raw device offers a more "direct" route to the physical device and allows an application more control over the timing of IO to that physical device. This makes raw devices suitable for complex applications like Database Management Systems that typically do their own caching.

Raw devices are character devices (major number 162). The first minor number (i.e. 0) is reserved as a control interface and is usually found at /dev/rawctl. A utility called **raw** (see **man raw**) can be used to bind a raw device to an existing block device. These "existing block devices" may be disks or cdroms/ dvds whose underlying interface can be anything supported by Linux (e.g. IDE/ATA or SCSI).

A sequence of commands listing the raw devices and then binding a SCSI disk partition followed by binding the whole disk looks like this on my system:

```
$ ls -lR /dev/raw*
crw-r--r--    1 root     root      162,   0 Dec  6 06:54 /dev/rawctl

/dev/raw:
total 0
crw-r--r--    1 root     root      162,   1 Dec  6 06:54 raw1
crw-r--r--    1 root     root      162,   2 Dec  6 06:54 raw2
crw-r--r--    1 root     root      162,   3 Dec  6 06:54 raw3
crw-r--r--    1 root     root      162,   4 Dec  6 06:54 raw4
$
$ raw -qa
$
$ raw /dev/raw/raw1 /dev/sda3
/dev/raw/raw1:  bound to major 8, minor 3
$ raw /dev/raw/raw2 /dev/sda
/dev/raw/raw2:  bound to major 8, minor 0
$ raw -qa
/dev/raw/raw1:  bound to major 8, minor 3
/dev/raw/raw2:  bound to major 8, minor 0
```

The normal array of system calls for character devices are available on raw devices. The size of the transfer for read(2) and write(2) must be an integral multiple of the physical device's block size. For a disk this will be its sector size which is normally 512 bytes. The data buffer given to read() and write() system calls must be aligned to the block size. The lseek(2) call needs to align its file read/write offset to a block boundary as well. The pread(3) call (see **man pread**) combines a read() and an lseek() and can be useful with raw devices (ditto with pwrite() ). Care should be taken with offsets greater than 2 GB (or perhaps 4 GB) on 32 bit architectures where the "off_t" type is 32 bits long. One solution is to use the _llseek() call (see **man llseek**).

Unix utilities such as recent versions of **dd** and **lmdd** (from the lmbench suite of programs) can be used to move data to and from "raw" devices as they meet the above-mentioned block alignment requirements. Recent versions of the **sg_dd** command in the sg_utils package can access both raw and sg devices.

Also note that if the physical device has an odd number of sectors (as shown by **blockdev --getsize /dev/ raw/raw***), the last sector will not be accessible using raw IO.

# Warning

If a block device is being accessed via a bound raw device and also via its normal block interface then there is no cache coherency between the two access mechanisms. For example if `/dev/sda1` was both mounted and being accessed via a bound raw device then there could be data inconsistencies.

# Chapter 12. Devfs pseudo file system

The main documentation for devfs can be found at: reference W5. The devfs name conventions for the SCSI subsystem are outlined in the section called "Device Names in devfs". Devfs is selected by the kernel build option CONFIG_DEVFS_FS and whether it is mounted at boot time (as `/dev`) or not is controlled by the kernel build option CONFIG_DEVFS_MOUNT. The latter option can be overridden by the kernel boot time options "devfs=mount" or "devfs=nomount", whichever is appropriate.

The devfs SCSI node names with their default permissions are:

```
disc          rw-------    whole disk including mbr
part1         rw-------    first partition {...p1}
...
part15        rw-------    15th partition {...p15}
cd            rw-rw-rw-    cd or dvd devices
mt            rw-rw-rw-    tape mode 0 with rewind {...m0}
mtl           rw-rw-rw-    tape mode 1 with rewind {...m1}
mtm           rw-rw-rw-    tape mode 2 with rewind {...m2}
mta           rw-rw-rw-    tape mode 3 with rewind {...m3}
mtn           rw-rw-rw-    tape mode 0 with no rewind {...m0n}
mtln          rw-rw-rw-    tape mode 1 with no rewind {...m1n}
mtmn          rw-rw-rw-    tape mode 2 with no rewind {...m2n}
mtan          rw-rw-rw-    tape mode 3 with no rewind {...m3n}
generic       rw-r-----
```

These node names are only present if the corresponding device (or sub-entities of the device (e.g. partitions)) and driver are present. For example if there is no sg driver present then there is no "generic" device name. The strings that appear above in braces are appended to the abridged "c0b0t0u0" notations outlined below as appropriate.

The devfs file names that are block or character special files will be called the primary device names in this description. The devfs daemon, called devfsd, introduces many symbolic links to those primary device names. This is done both for backward compatibility and convenience. These symbolic links will be called secondary device names.

The secondary device names are controlled by the devfsd configuration file usually found in `/etc/devfsd.conf`. Following is a list of secondary device names when the default devfsd.conf file is used:

```
Secondary name         slink to this primary device name
-------------------------------------------------------
/dev/sda               /dev/scsi/host0/bus0/target2/lun0/disc
/dev/sda1              /dev/scsi/host0/bus0/target2/lun0/part1
/dev/sd/c0b0t2u0       /dev/scsi/host0/bus0/target2/lun0/disc
/dev/sd/c0b0t2u0p1     /dev/scsi/host0/bus0/target2/lun0/part1
/dev/sr0               /dev/scsi/host0/bus0/target4/lun0/cd
/dev/sr/c0b0t4u0       /dev/scsi/host0/bus0/target4/lun0/cd
/dev/st0               /dev/scsi/host1/bus0/target0/lun0/mt
/dev/nst0a             /dev/scsi/host1/bus0/target0/lun0/mtan
/dev/st/c1b0t0u0m0     /dev/scsi/host1/bus0/target0/lun0/mt
/dev/st/c1b0t0u0m3n    /dev/scsi/host1/bus0/target0/lun0/mtan
/dev/sg0               /dev/scsi/host0/bus0/target2/lun0/generic
```

```
/dev/sg1               /dev/scsi/host0/bus0/target4/lun0/generic
/dev/sg2               /dev/scsi/host1/bus0/target0/lun0/generic
/dev/sg/c0b0t2u0       /dev/scsi/host0/bus0/target2/lun0/generic
/dev/sg/c0b0t4u0       /dev/scsi/host0/bus0/target4/lun0/generic
/dev/sg/c1b0t0u0       /dev/scsi/host1/bus0/target0/lun0/generic
```

Note that the more common `/dev/scd0` variant for SCSI cdroms is not supported. There are also `/dev/discs`, `/dev/cdroms` and `/dev/tapes` directories that contain symbolic links to all devices (i.e. not just SCSI devices) that fall into that categorization:

```
Secondary name        slink to this primary device
---------------------------------------------------------
/dev/discs/disc0      /dev/ide/host0/bus0/target0/lun0      *
/dev/discs/disc1      /dev/scsi/host0/bus0/target2/lun0     *
/dev/cdroms/cdrom0    /dev/ide/host0/bus1/target1/lun0/cd
/dev/cdroms/cdrom1    /dev/scsi/host0/bus0/target4/lun0/cd
/dev/tapes/tape0      /dev/scsi/host1/bus0/target0/lun0     *
```

Those entries marked with "*" are directories containing the primary devices. Note that IDE/ATA devices are listed before SCSI devices. These secondary device names mimic the same persistence rules as the primary device names. So when a SCSI device (?), or its lower level driver or its upper level driver are removed then so are the primary and secondary device names associated with it.

When devfs is mounted as `/dev`, the old "`/dev/sda6`" type can still be used in some contexts. This may be convenient if typing is required at the kernel boot time prompt. For example if a user wants to change the root partition on a "devfs" machine then any of the following examples may be used as a kernel boot time option:

```
root=/dev/sda6
root=/dev/scsi/host0/bus0/target0/lun0/part6
root=/dev/sd/c0b0t0u0p6
```

There are many device scanning programs that expect to see the pre-devfs device names present and it will some time before they become devfs aware. Also some programs rely on a open of `/dev/sg0` (for example) to load the sg driver (assuming it is a module and not already loaded). This can be arranged by an entry in `/etc/devfsd.conf` file of:

```
        LOOKUP          sg.*            MODLOAD
```

and the following in `/etc/modules.devfs`:

```
        probeall        /dev/sg         scsi-hosts sg
        alias           /dev/sg*        /dev/sg
```

The sg device permissions can be changed with this entry in the `/etc/devfsd.conf` file:

```
  REGISTER scsi/host.*/bus.*/target.*/lun.*/generic
                        PERMISSIONS 0.0 rw-rw-rw-
```

See "man devfsd" for more information.

An application can determine whether devfs is active by the presence or otherwise of the file `/dev/.devfsd`.

A feature of a /dev directory based on a persistent file system (e.g. ext2) is the ability to associate permissions with a device file name and keep them from one boot to the next. As noted above the default action of devfs is to assign device file name permissions anew each time a machine is booted. The PERMISSIONS action in the `/etc/devfsd.conf` can be used to assert permissions but this may be considered a little awkward. The devfs document (W5) describes a method for getting the best of both worlds. This technique relies on the recently added feature in lk 2.4 to mount the same file system at multiple points.

# Appendix A. Common bus types (SCSI and other)

A very good overview of the various bus types touched on in this appendix (both SCSI and others) can be found at `www.pctechguide.com/04disk2.htm` [http://www.pctechguide.com/04disk2.htm].

**SCSI.** The original SCSI 1 standard (ANSI specification X3.131-1986) introduced an 8 bit parallel bus that was able to do asynchronous transfers at 1.5 MegaBytes/sec and synchronous transfers up to 5 MB/sec. SCSI commands are sent at the asynchronous rate. SCSI data is transferred either at the asynchronous rate (worst case) or a negotiated synchronous rate (with 5 MB/sec being the best case).

**FAST SCSI.** The SCSI 2 standard raised the maximum synchronous speed to 10 MB/sec. SCSI 2 defined several parallel buses: single ended (as used by SCSI 1) and a new differential bus. The differential bus has better noise immunity and its maximum bus length is 25 metres (compared with single ended's 6 metres). Tagged queuing of commands was also added by SCSI 2.

**WIDE SCSI.** The SCSI 2 standard also increased the width of the bus allowing 16 and 32 bit "wide" variants. Very little use has been made of the 32 bit width so "wide" usually refers to a 16 bit wide data path. The maximum number of SCSI devices that can connect to a parallel SCSI bus is directly related to the bus width hence "wide" buses allow a maximum of 16 SCSI devices to be connected. [At least one of those devices must be the SCSI "initiator" which is usually a host adapter.]

**ULTRA SCSI.** Traditionally synchronous buses are clocked either on the rising or falling edge of the clock (which is normally a square wave). A recent trend has been to clock on both edges and thus double the available bandwidth. This is how ULTRA SCSI doubles the SCSI 2 "fast" speed to 20 MB/sec.

**ULTRA WIDE SCSI.** The same "ultra" technique applied to a (16 bit) wide SCSI parallel bus yields a bandwidth of 40 MB/sec.

**ULTRA 2 WIDE SCSI.** This variant introduces a new "low voltage" differential signalling (LVD) that allows the synchronous clock speed to be doubled yielding 80 MB/sec when using a (16 bit) wide bus. In this case the maximum SCSI bus length is 12 metres. To be backward compatible with ULTRA WIDE this variant can fall back to "single ended" operation. This leads to the abbreviation LVD/SE being used by adapter manufacturers. One shortcoming of this approach is that the presence of one UW device on a U2W bus will cause all other U2W devices to communicate at the slower (i.e. UW) rate. Some adapters overcome this by having separate LVD and SE physical buses on the same logical SCSI bus.

**ULTRA 160 SCSI.** ULTRA 160 doubles parallel SCSI bus bandwidth yet again. It uses a 16 bit wide data path, LVD signalling (see previous entry) and double transition clocking that increases the maximum synchronous bandwidth to 160 MB/sec. Additional features include cyclic redundancy codes (CRC) to improve data integrity (compared with a parity bit) and domain validation which adjusts transfer rates if the error rate is too high.

**ULTRA 320 SCSI.** Shortly ULTRA 320 adapters will be available (disks with that interface are already on the market). This is also a 16 bit wide LVD bus that can fall back to slower speeds for compatibility with older devices. It extends the features of Ultra 160 by doubling the clock speed. Packetized SCSI which sends commands and status at full bus speed (rather than 5 MB/sec) is included. Other improvements include "quick arbitration and selection", "read and write data streaming" and CRC protection for command blocks as well as data (Ultra 160 had CRC protection for data only). Note that adapter cards using 64 bit PCI (or better: PCI-X) are required to stop the PCI bus being a bottleneck at these speeds. More information can be found at `www.scsita.org` [http://www.scsita.org]. Recently an Ultra 320 HBA vendor claimed up to 105,000 IO operations per second which implies per command SCSI bus overhead is less than 10 microseconds. There is a draft ULTRA 640 standard but that may be overtaken by Serial Attached SCSI.

**Serial Attached SCSI (SAS).**    Serial Attached SCSI (SAS) uses the same transport technology as Serial ATA (sATA) and extends it. [sATA is described below.] SAS can use external expanders to control up to 16000 devices from a single HBA. The data transfer is full duplex and 1.5 Gbps or 3.0 Gbps "phys"s can be aggregated to increase bandwidth. Cable lengths can be up to 6 meters. SAS disks are dual ported. sATA disks can be connected to a SAS expanders (but SAS disks can't be connected to a sATA HBA). SAS was demonstrated at CeBit recently but won't reach the market until 2004.

**FC-AL.**    This stands for Fibre Channel - Arbitrated Loop and may involve dual 2 Gigabit per second single mode fibre optic links spanning 10 kilometres with throughput of up to 400 MegaBytes per second. Often associated with storage area networks (SANs). Up to 126 devices can be attached to a loop which in turn can be extended to 16 million devices in public loop mode. The transmission medium isn't necessarily fibre optic cable: copper (in the form of co-axial cable) can also be used at lower speeds and for shorter distances.

**SRP/InfiniBand.**    SRP (SCSI RDMA Protocol) [ `SRP_draft` [http://www.t10.org/drafts.htm#SCSI3_SRP]] is a SCSI transport for InfiniBand [ `Infiniband_trade_association` [http://www.infinibandta.org]], a high-performance interconnect running at 10 and 30 Gbps. SRP driver source is available at `infiniband.sourceforge.net/Storage/SRP` [http://infiniband.sourceforge.net/Storage/SRP/index.htm] .

**IEEE 1394.**    This standard also goes by the name of "Fire Wire" [trademarked by Apple] and "iLink" [trademarked by Sony]. It is a serial bus that can run at up to 400 Megabits/sec (IEEE 1394a). A newer standard, IEEE 1394b, ups this to 800 Megabits/sec (with extensions to 1.6 and 3.2 Gigabits/sec) with cable runs up to 100 metres. It has a similar but more general architecture than USB. The IEEE 1394 standard allows for the SCSI command set to be carried over a 1394 bus. There is a "sbp2" driver now available for the Linux IEEE 1394 stack. This sbp2 driver is also a SCSI subsystem lower level driver (so it is functionally similar to the ide-scsi driver). So IEEE 1394 devices that use the SBP-2 protocol (e.g. disks, cd-rw/dvd drives, MO drives and scanners) can be accessed via the SCSI subsystem. See `Linux1394.sourceforge.net` [http://Linux1394.sourceforge.net] for more information. The sbp2 driver is now in lk 2.4.7 .

**iSCSI.**    This is a new IETF standard for sending the SCSI command set over a TCP connection (or several of them). This permits SCSI devices (targets such as disks) to be network appliances, accessed locally (or potentially at a great distance) by a host machine.

**NON SCSI buses.**    The following buses are not defined by the SCSI standards but are of interest because they either can carry the SCSI command set, are in some way related to the Linux SCSI subsystem or supply a similar functionality to SCSI products.

**IDE/ATA (ATAPI).**    IDE is the most used disk type on PC systems today. The acronym stands for Integrated Drive Electronics and as the name suggests it places the bulk of the IO "intelligence" on the disk controller card rather than spreading it between the device (most often a disk) and a controller (HBA) as SCSI does. IDE grew out of the ST506 and ESDI standards in the 1980s. EIDE (extended IDE) is a related acronym. The modern standards that refer to this bus architecture are known as ATA and can be found at `www.t13.org` [http://www.t13.org]. The ATA Packet Interface (ATAPI) extends the disk oriented command set to support CDROM and tape drives. The ATAPI command set closely resembles the SCSI command set. The most recent ATA technology is outlined in the next paragraph.

**ATA 133.**    The ATA standards used by IDE devices have also been marching through the adjectives (e.g. fast and ultra) and the numbers (e.g. 2, 33, 66, 100 and 133). The most recent addition is ATA 133 which supports burst rates of 133 MB/sec and up to 2 devices per bus. [PCs typically have 2 and often 4 ATA buses.] ATA 66, 100 and 133 need a special cable. ATA cables are relatively short precluding IDE devices being external to the computer. Cable lengths have previously been limited to 18 inches although 1 metre long cables have now appeared. Coincidently 133 MB/sec in also the maximum throughput of the normal PCI bus found in most PCs. The are higher speed (and wider) versions of PCI but they are relatively rare.

**Serial ATA (sATA).**     Serial ATA uses 2 differential pairs to exchange data with a sATA disk less than 1 metre away at 1.5 Gigabits per second. One pair takes data to the disk and the other returns data from the disk. Data rates up to 150 Megabytes per second are possible (data transfer is half duplex). sATA is a point to point connection, not a bus, so ATA's master and slave strapping disappears. sATA cabling is less bulky and the form factor of its plugs and sockets are smaller than parallel ATA (and the SCSI Parallel interface). sATA devices are beginning to appear on the market. sATA-2 is a draft standard that doubles the serial data rate to 3 Gigabits per second.

**USB.**     Universal Serial Bus (USB) has a bandwidth of between 1.5 and 12 Megabits/sec (the latter speed with USB 1.1). Up to 127 devices can be connected using a series of hubs each of which connects up to 7 devices (with a 5 metre limit). USB supplies 5 volts at 0.5 amps to power small devices. USB is "plug and play", hot pluggable and supports isochronous data transfers (required for audio and video devices that need guaranteed minimum bandwidth).

**PC Parallel port.**     The original PC parallel port was uni-directional (towards the printer) and was capable of about 10 KB/sec. The IEEE 1284 standard in 1994 introduced 5 modes of data transfer:

- Compatibility mode (forward direction)

- Nibble mode (reverse direction)

- Byte mode (reverse direction)

- EPP mode (bi-directional)

- ECP mode (bi-directional)

Enhanced Parallel Port (EPP) achieves transfer speeds of between 500 KB/sec and 2 MB/sec and is targeted at CD-ROMs, tapes and hard drives. Extended Capability Port (ECP) includes run length encoding and support for DMA. ECP is targeted at fast printers and scanners.

**I2O.**     "The I2O (Intelligent Input/Output) specification defines a standard architecture for intelligent I/O that is independent of both the specific device being controlled and the host operating system (OS)" [from `www.i2osig.org` [http://www.i2osig.org]]. It defines a "split driver" model in which the OS Services Module (OSM) sits between the host OS device interface and the I2O communications layer while the Hardware Device Module (HDM) sits between the I2O communications layer and the hardware. The HDM may well run on a dedicated processor (IOP).

# Appendix B. Changes between lk 2.2 and (during) 2.4

Significant work has been done to change the single SCSI command queue used in lk 2.2 to one command queue per device. To make the SCSI subsystem more SMP friendly the granularity of the locks is much finer grained. In lk 2.2 the whole subsystem essentially used one lock.

Even though it is not part of the SCSI subsystem, the inclusion of devfs solves many SCSI device addressing problems that existed in the past. Associated with devfs but very useful even in its absence is the "scsihosts" kernel boot time (and module load time) option. This option allows users to have some control over the ordering of multiple SCSI hosts.

This appendix is difficult to maintain since features and drivers that have proven useful in lk 2.4 (and its development tree) have tended to be back ported into the higher release numbers of the lk 2.2 series.

Currently (lk 2.4.2) support for MO devices is broken. Old DOS file systems with a block size of 2048 bytes also have been reported as broken. The problem seems to arise with media that have a physical block size larger than the 1 KB logical block size used by the block subsystem. Only the sd driver has this problem (luckily not the sr driver in which 2048 byte sectors are the norm).

## Mid level changes

```
SCSI_IOCTL_GET_IDLUN      {ioctl, changed}
```

## sd changes

```
HDIO_GETGEO_BIG           {ioctl, new}
```

## sr changes

No sr changes reported. As a related matter, the "hdx=scsi" kernel boot option has been added. See the section called "ATAPI cdroms" for more details.

## st changes

No interface changes. In lk 2.2 the maximum number of extra tape devices that could be added after boot time was limited to 3. This limitation has been removed (leaving a maximum of 32 tape devices as noted earlier).

A variant st driver called "osst" to handle early model OnStream tape drives has been added in lk 2.4 .

## sg changes

The main change is the addition of a new interface structure called "sg_io_hdr". The existing interface structure (called "sg_header") was found to be inflexible requiring the concatenation of raw data together with meta-data in the read() and write() commands.

```
sg_io_hdr              {new interface structure}
SG_IO                  {new ioctl}
direct IO              {present but commented out, see ALLOW_DIO}
procfs output          {new information in /proc/scsi/sg directory}
boot/module parameters  {new}
```

Up to 64 bytes of sense data can be obtained from the sg_io_hdr interface structure. Also a residual count associated with the data transfer is available (if the lower level driver supports it, if not the residual count will be 0).

# Changes during the lk 2.4 series

Even though the lk 2.4 production series is meant to be "stable" there have been a significant number of changes as well as bug fixes. The following list does not include changes to the lower level (adapter) drivers. Each item of the list is prefixed by the kernel version that it was introduced. [1]

- [2.4.4] added the SCSI_IOCTL_GET_PCI ioctl(),

- [2.4.7] the "lun" bits (3 bits representing lun values 0 through 7 in the SCSI 1 and SCSI 2 standards) are no longer masked into the second byte of SCSI commands if the INQUIRY for that devices shows a SCSI level greater than SCSI_2,

- [2.4.7] the max_scsi_luns kernel (and module scsi_mod) option previously could be 1 to 7. Now the upper value can be large. [The scan algorithms are still doing a sequential scan rather than using REPORT_LUNS.]

- [2.4.7] both scsi_unregister_host() and scsi_unregister_module() now return an int (previously they were void functions). They return 0 for success, -1 for failure (typically busy),

- [2.4.7] the upper level drivers now report the correct scsi device name when they are attached. [The log messages that started with "Detected ..." previously sometimes reported the wrong device (e.g. sdc rather than sdb).] Kernel boot up messages will now show SCSI devices as "Attached ...",

- [2.4.7] 'max_sectors' was added to the Scsi_Host structure,

- [2.4.8] some mid level logic was altered to retry commands if the sense buffer indicates that logical unit is becoming ready [ASC=4, ASQ=1],

- [2.4.9] a major st update,

- [2.4.9] mid level changed to retry commands if lower level (adapter) driver returned DID_RESET,

- [2.4.10] original result (including SCSI status) saved when mid level issues a REQUEST SENSE so it can be restored afterwards,

- [2.4.10] added BLKGETSIZE64, BLKBSZSET and BLKBSZGET ioctls to sd + sr,

- [2.4.10] sg update that fixes generic_unplug_device() race + bumps access_count on opens (and decrements on releases),

---

[1] This list has been compiled from the official 2.4 series kernels released at  www.kernel.org [http://www.kernel.org]. Distributions are free to tailor the official kernels and this may impact what is supported (or changed) in the SCSI subsystem. For example this machine reports this kernel: "2.4.18-27.8.0". So that is roughly based on the official 2.4.18 kernel which the vendor has "modded" 27 times for the "8.0" level of their distribution. As an example of the type of changes, the aic7xxx driver in the official 2.4.18 does not support Adaptec's Ultra 320 series of PCI adapters; however that vendor's version does.

- [2.4.11] added MODULE_LICENSE macro in most drivers, mostly MODULE_LICENSE("GPL"),

- [2.4.11] scsi_pid bumped for each command (why?),

- [2.4.11] st update to bump access_count. Now all upper level drivers increment access_count on opens and decrement it on releases,

- [2.4.13] scatterlist structure grows (alt_address is removed, page and offset added),

- [2.4.13] don't probe luns > 7 for target <= SCSI_2 ,

- [2.4.14] fine tuning (bug fixes) associated with scatterlist structure changes [it broke st ?],

- [2.4.15] 16 byte SCSI commands permitted [MAX_COMMAND_SIZE changes from 12 to 16]. HBA driver must set Scsi_Host::max_cmd_len to 16 for mid level to forward 16 byte SCSI commands,

- [2.4.15] BLKGETSIZE + BLKGETSIZE64 ioctl() implementations moved out of SCSI subsystem (and into block subsystem),

- [2.4.15] large st update,

- [2.4.15] lk 2.5.0 forks off so lk2.4.15==lk2.5.0 .

- [2.4.17] add generic_unplug_device() call to scsi_wait_req(). This stops long waits in SCSI_IOCTL_SEND_COMMAND.

- [2.4.17] fix device scanning bug where, in some cases, the scsi_level (i.e. SCSI standard adherence) was misplaced.

- [2.4.17] major sg driver update, add mmap()-ed IO

- [2.4.18] permit upper level driver "init()" functions (e.g. sd_init() ) to fail gracefully. [Add Scsi_Device::detected and scsi_unregister_module() .]

- [2.4.18] Fix for clustering (SCSI commands) on MO devices.

- [2.4.18] st driver update (compression algorithms).

- [2.4.18] update Documentation/scsi.txt and scsi-generic.txt .

- [2.4.18] Revamp scsi_debug driver .

- [2.4.19] Scsi reservation and reset capability added. Reservations allow multiple machines to share the same device (via a reserve/release mechanism). Scsi reset (via sg) is needed to "break" a reservation held by a non-responding machine.

- [2.4.19] Introduce BLIST_LARGELUN to handle LUNs larger than 7 despite reporting SCSI 2.

- [2.4.19] Change sd and sr so RECOVERED_ERROR is not treated as a hard error. Send warning to log/console.

- [2.4.19] Zero out sg's buffers before use. [Sg version upgraded from 3.1.22 to 3.1.24 but this is not reflected in sg.h (supeficial).]

- [2.4.20] Support for highmem I/O added. Used by aic7xxx, 3w-xxxx, esp, megaraid, qlogicfc and sym53c8xx_2 LLDs.

- [2.4.20] "blocking_open" boot time, module load time parameter added to st.

- [2.4.21] Give new HBAs a new host number (higher than any previously used) unless there is a "scsi-hosts" match. Host numbering sequence "holes" are only re-used if there is a "scsihosts" match.

- [2.4.21] stop the SCSI status RECOVERED ERROR being treated as an error by the mid level (complements a change in 2.4.19).

- [2.4.21] use the TEST_UNIT_READY command (rather than START_STOP) to determine if removable media has changed (in sd driver).

- [2.4.21] major work on ide-scsi driver.

- [2.4.21] add aic79xx driver for Adaptec Ultra 320 controllers.

- [2.4.22] Extend timeout of SEND DIAGNOSTIC command to 2 hours. This is for foreground extended self tests.

- [2.4.26] Add 'scsi_allow_ghost_devices' kernel boot time and scsi_mod module option.

- [2.4.27] SATA support via "libata" library introduced. SATA disks appear with SCSI subsystem names (e.g. "/dev/sdb) and respond to SCSI commands (via a command translation facility).

# Appendix C. Troubleshooting

Many SCSI problems are caused by cabling and (lack of, or inappropriate) termination. This often results in repeated SCSI bus resets, parity or CRC errors and sometimes reduced transfer speeds. There is a good SCSI termination tutorial at this site: www.scsita.org/aboutscsi/SCSI_Termination_Tutorial.html [http://www.scsita.org/aboutscsi/SCSI_Termination_Tutorial.html]. There is other useful SCSI information at that site (see W9).

There is also a SCSI "faq" site (see W10) that addresses many configuration and troubleshooting issues. Although the main focus of this site is Windows (and its ASPI interface), much is relevant to SCSI in Linux and other Unix implementations.

When it looks like something has partially locked up the system, the **ps** command can be useful for finding out what may be causing the problem. The following options may be useful for identifying what part of the kernel may be causing the problem. This information could be forwarded to the maintainers.

```
ps -eo cmd,wchan
ps -eo fname,tty,pid,stat,pcpu,wchan
ps -eo pid,stat,pcpu,nwchan,wchan=WIDE-WCHAN-COLUMN -o args
```

The most interesting option for finding the location of the "hang" is "wchan". If this is a kernel address then **ps** will use `/proc/ksyms` to find the nearest symbolic location. The "nwchan" option outputs the numerical address of the "hang".

If the system is not responding to keystrokes, then <Alt+ScrollLock> in text mode should output a stack trace while <Ctrl+ScrollLock> should output a list of all processes. If the log is still working, the output will be sent there as well as appearing on the console.

If the kernel has been built with the CONFIG_MAGIC_SYSRQ, then in text mode <Alt+SysRq+H> will list available commands. Of these <Alt+SysRq+S> is useful for doing an emergency sync while <Alt+SysRq+U> will remount file systems in read only mode. After that <Alt+SysRq+B> to reboot the machine might be your next move.

# Appendix D. Performance, Test and Debugging tools

**scu.**      The SCSI Command Utility (SCU) implements various SCSI commands necessary for normal maintenance and diagnostics of SCSI peripherals. Some of its features include: formatting, scanning for (and reassigning) bad blocks, downloading new firmware, executing diagnostics and obtaining performance information. It is available on several Unix platforms (and NT), however it is only currently available in binary form. See  www.bit-net.com/~rmiller/scu.html [http://www.bit-net.com/~rmiller/scu.html] for more details.

**dd.**      Very useful for testing the streaming performance of disks and cdroms/dvds. See **man dd** for more details. Here is an example for timing how long a disk takes to read 1 GB (10**9 bytes) starting from block 0:

```
$ time dd if=/dev/sda of=/dev/null bs=512 count=1953126
```

If the raw device `/dev/raw/raw1` is bound to `/dev/sda` then the above line is equivalent to:

```
$ time dd if=/dev/raw/raw1 of=/dev/null bs=512 count=1953126
```

This may be slower than expected since one 512 byte sector is being read at a time. Changing the last 2 arguments to "bs=8k count=122071" should give better timings for the "raw" dd.

**dt.**      The Data Test (DT) program is modelled on dd's syntax but dt can do a lot more than sequential copies. It is a comprehensive data test program for SCSI devices such as disks, tapes and cdrom/dvds. It is available on several Unix platforms (and NT), and its source is available (unlike its stable mate "scu" discussed earlier). See  www.bit-net.com/~rmiller/dt.html [http://www.bit-net.com/~rmiller/dt.html] for more details.

**lmdd.**      This command is part of the lmbench suite of programs and is a variant of the **dd** command. It has been tailored for IO measurements and outputs timing and throughput numbers on completion. Hence the **time** command and a calculator are not needed.

**blockdev.**      Fetches the sector size, the number of sectors and read ahead status of a block device (typically a disk). Can also be used to flush buffers and reread the partition table. See **man blockdev**.

**sg_dd.**      This command is part of the sg_utils package (see W4) and is another variant of the **dd** command in which either the input and/or output file is a sg or a raw device. The block size argument ("bs") must match that of the physical device in question. The "skip" and "seek" arguments can be up to 2**31 - 1 on a 32 bit architecture allowing 1TB disks to be accessed (2G * 512). The Linux system command llseek() is used to seek with a 64 bit file read/write offset. The **lmdd** does not handle the > 2GB case and the **dd** command gets creative with multiple relative seeks. **sg_dd** has a "bpt" (blocks per transfer) argument that controls the number of blocks read or written in each IO transaction.

There are other programs in the sg_utils package to scan the SCSI bus (**sg_scan** and **sg_map**), to measure SCSI bus throughput (**sg_rbuf** and **sg_turs** ), show data from the SCSI inquiry command (**sg_inq**) and spin up (or down) media (**sg_start**).

**dd_rescue + scsiinfo.**      This dd variant is designed to rescue damaged media such as SCSI (or IDE) disks and CDROMs (see W6). The  **scsiinfo** utility for displaying and changing mode page information is also at that site.

**sard.**     This utility is modelled on System V Release 4's **sar -d** for producing IO statistics for mounted devices and partitions. It has been developed by Stephen Tweedie and includes the sard utility and a required kernel patch which expands the output of `/proc/partitions`. It can be found at ftp.uk.linux.org/pub/linux/sct/fs/profiling [ftp://ftp.uk.linux.org/pub/linux/sct/fs/profiling]. It collects statistics at a relatively low level (e.g. SCSI mid level) compared to programs like **vmstat** (see "man vmstat").

# Appendix E. Compile options and System calls including ioctls

The compile options in this appendix are those which a system administrator might conceivably want to change. Naturally the defaults are chosen so the vast majority of users will not need to modify anything. In some cases setting kernel build time options, kernel boot time parameters or module load parameters has the same effect as changing a driver compile time option.

System calls act as the interface between application programs and the kernel and its drivers. In the case of the layered driver architecture that the SCSI subsystem uses, the upper layer drivers handle most of the system calls.

The SCSI subsystem has a "bubble down" ioctl structure. First the upper level driver associated with the open file descriptor attempts to decode the ioctl. If it doesn't recognize it then the ioctl is passed down to the mid level. If the mid level doesn't recognize it then the ioctl is passed down to the lower level driver associated with the file descriptor. If the lower level driver doesn't recognize it then a EINVAL error is generated.

Some ioctls are dispatched to related subsystems.

# Mid level

The following header files in the kernel source are relevant to the mid level:

```
/usr/src/linux/include/scsi/scsi.h
/usr/src/linux/include/scsi/scsi_ioctl.h
```

These files are meant for applications to use (other than parts in a __KERNEL__ conditional compilation block). They may also be found in /usr/include/scsi directory but it is best not to trust these versions as they are maintained with the glibc library and may lag the kernel version being used. Usually in Linux systems `/usr/include/linux` can be relied upon to be a symbolic link to the kernel source's include area (typically `/usr/src/linux/include/linux`). This symbolic link can be used to include the correct `scsi_ioctl.h` using the following trick: `#include <linux/../scsi/scsi_ioctl.h>`

This include file: `/usr/src/linux/drivers/scsi/scsi.h` is the key internal header file for the SCSI subsystem. As such it will not be discussed here other than to point out it has the same file name (but it's in a different directory) as the include file mentioned at the beginning of this section. This sometimes causes confusion.

The mid level `drivers/scsi/scsi_scan.c` file maintains an array of known SCSI devices with *idiosyncrasies* . [This was known as the "black list" but that was considered to judgmental.] The array is called "device_list". The various value are:

- BLIST_NOLUN only probe lun 0

- BLIST_FORCELUN force all 8 luns to be probed

- BLIST_BORKEN passes through broken flag to lower level driver

- BLIST_KEY sends magical MODE SENSE (pc=0x2e) to unlock device

- BLIST_SINGLELUN only allow IO on one lun at a time

- BLIST_NOTQ disable tagged queuing

- BLIST_SPARSELUN keep going after lun not found

- BLIST_MAX5LUN only probe up to lun 5

- BLIST_ISDISK override INQUIRY's type with disk (direct access) type

- BLIST_ISROM override INQUIRY's type with ROM

# Mid level compile options

None.

# Mid level ioctls

See the following files:

```
/usr/src/linux/include/scsi/scsi.h
```

Note that the SCSI status constants defined in include/scsi/scsi.h are shifted 1 bit right from the values in the SCSI standards:

```
scsi.h constant        value      SCSI 2 standard value
-----------------------------------------------------
CHECK_CONDITION        0x1          0x2
CHECK_GOOD             0x2          0x4
BUSY                   0x4          0x8
....
```

Summary of ioctl()s follow:

```
SCSI_IOCTL_SEND_COMMAND
  This interface is deprecated - users should use
  the scsi generic (sg) interface instead, as this
  is a more flexible approach to performing
  generic SCSI commands on a device.

  The structure that we are passed should look like:

  struct sdata {
   unsigned int inlen;      [i] Length of data written to device
   unsigned int outlen;     [i] Length of data read from device
   unsigned char cmd[x];    [i] SCSI command (6 <= x <= 16)
                            [o] Data read from device starts here
                            [o] On error, sense buffer starts here
   unsigned char wdata[y]; [i] Data written to device starts here
  };
  Notes:
    -  The SCSI command length is determined by examining
```

the 1st byte of the given command. There is no way
to override this.
- Data transfers are limited to PAGE_SIZE (4K on
  i386, 8K on alpha).
- The length (x + y) must be at least OMAX_SB_LEN
  bytes long to accommodate the sense buffer when
  an error occurs. The sense buffer is truncated to
  OMAX_SB_LEN (16) bytes so that old code will not
  be surprised.
- If a Unix error occurs (e.g. ENOMEM) then the user
  will receive a negative return and the Unix error
  code in 'errno'. If the SCSI command succeeds then
  0 is returned. Positive numbers returned are the
  compacted SCSI error codes (4 bytes in one int)
  where the lowest byte is the SCSI status. See the
  drivers/scsi/scsi.h file for more information on this.

SCSI_IOCTL_GET_IDLUN
  This ioctl takes a pointer to a "struct scsi_idlun" object
  as its third argument. The "struct scsi_idlun" definition
  is found in <scsi/scsi.h>. It gets populated with scsi
  host, channel, device id and lun data for the given device.
  Unfortunately that header file "hides" that structure
  behind a "#ifdef __KERNEL__" block. To use this, that
  structure needs to be replicated in the user's program.
  Something like:
  typedef struct my_scsi_idlun {
       int four_in_one;     /* 4 separate bytes of info
                               compacted into 1 int */
       int host_unique_id; /* distinguishes adapter cards from
                               same supplier */
  } My_scsi_idlun;
       "four_in_one" is made up as follows:
       (scsi_device_id | (lun << 8) | (channel << 16) |
       (host << 24))
  These 4 components are assumed (or masked) to be 1 byte each.

SCSI_IOCTL_GET_BUS_NUMBER
  In lk 2.2 and earlier this ioctl was needed to get the
  host number. During lk 2.3 development the
  SCSI_IOCTL_GET_IDLUN ioctl was changed to include this
  information. Hence this ioctl is only needed for
  backward compatibility.
SCSI_IOCTL_TAGGED_ENABLE
  Probably a remnant of the past when the mid level
  addressed such issues. Now this functionality is
  controlled by the lower level drivers. Best ignored.
SCSI_IOCTL_TAGGED_DISABLE
  See comment for SCSI_IOCTL_TAGGED_ENABLE.
SCSI_IOCTL_PROBE_HOST
  This ioctl expects its 3rd argument to be a pointer to
  a union that looks like this:
  union probe_host {
    unsigned int length;  /* [i] max length of

```
                                    output ASCII string */
  char str[length];      /* [o] N.B. may need '\0'
        appended */
};
The host associated with the device's fd either has a
host dependent information string or failing that its
name, output into the given structure. Note that the
output starts at the beginning of given structure
(overwriting the input length). N.B. A trailing '\0'
may need to be put on the output string if it has been
truncated by the input length. A return value of 1
indicates the host is present, 0 indicates that the
host isn't present (how can that happen?) and a
negative value indicates an error.
```

```
SCSI_IOCTL_DOORLOCK
SCSI_IOCTL_DOORUNLOCK
SCSI_IOCTL_TEST_UNIT_READY
  Returns 0 if the unit (device) is ready, a positive
  number if it is not or a negative number when there
  is an OS error.
```

```
SCSI_IOCTL_START_UNIT
SCSI_IOCTL_STOP_UNIT
SCSI_EMULATED_HOST             {same as SG_EMULATED_HOST <new>}
```

```
SCSI_IOCTL_GET_PCI
  Yields the PCI slot name (pci_dev::slot_name) associated with the lower
  level (adapter) driver that controls the current device. Up to 8 characters
  are output to the locations pointed to by 'arg'. If the current device
  is not controlled by a PCI device then errno is set to ENXIO.
  [This ioctl() was introduced in lk 2.4.4]
```

# sd driver

## sd compile options

```
MAX_RETRIES     {5}
SD_TIMEOUT      {30 seconds}
SD_MOD_TIMEOUT  {75 seconds}
```

## sd ioctls and user interface

The relevant files to see:

```
include/linux/hdreg.h
include/linux/genhd.h
include/linux/fs.h
```

A list of ioctl()s follow:

```
HDIO_GETGEO_BIG
HDIO_GETGEO      [retrieve disk geometry]
BLKGETSIZE       [number of sectors in device]
BLKROSET         [set read only flag]
BLKROGET         [get read only flag]
BLKRASET         [set read ahead value]
BLKRAGET         [get read ahead value]
BLKFLSBUF        [instructs SCSI subsystem to flush buffers]
BLKSSZGET        [get device block size]
BLKPG   [partition table manipulation]
BLKELVGET        [get elevator parameters]
BLKELVSET        [set elevator parameters]
BLKRRPART        [reread the partition table]

open()    (all flags ignored)
close()
ioctl()   (see list above)
```

# sr driver

## sr compile options

None.

## sr ioctls and user interface

See the following files:

```
/usr/src/linux/include/linux/cdrom.h
/usr/src/linux/drivers/cdrom/cdrom.c [revision history section]
/usr/src/linux/Documentation/cdrom/cdrom-standard.tex
```

Some of the following ioctls are described in cdrom-standard.tex :

```
CDROMCLOSETRAY
CDROM_SET_OPTIONS
CDROM_CLEAR_OPTIONS
CDROM_SELECT_SPEED
CDROM_SELECT_DISC
CDROM_MEDIA_CHANGED
CDROM_DRIVE_STATUS
CDROM_CHANGER_NSLOTS
CDROM_LOCKDOOR
CDROM_DEBUG
CDROM_GET_CAPABILITY
DVD_READ_STRUCT
DVD_WRITE_STRUCT
DVD_AUTH
CDROM_SEND_PACKET
```

```
CDROM_NEXT_WRITABLE
CDROM_LAST_WRITTEN
```

The O_NONBLOCK flag on the open() of scd devices is important. Without it the open() will wait until there is media in the device before returning.

```
open()              O_NONBLOCK
close()
read()
write()
ioctl()
```

# st driver

## st compile options

Most of the following compile options can be overridden with boot/module parameters and/or runtime configuration (i.e. ioctls).

The following parameters are defined in linux/drivers/scsi/st_options.h

```
ST_NOWAIT                     {0}
ST_IN_FILE_POS                {0}
ST_RECOVERED_WRITE_FATAL      {0}
ST_DEFAULT_BLOCK              {0}
ST_BUFFER_BLOCKS              {32}
ST_WRITE_THRESHOLD_BLOCKS     {30}
ST_MAX_BUFFERS                {4}
ST_MAX_SG                     {16}
ST_FIRST_SG                   {8}
ST_FIRST_ORDER                {5}
ST_TWO_FM                     {0}
ST_BUFFER_WRITES              {1}
ST_ASYNC_WRITES               {1}
ST_READ_AHEAD                 {1}
ST_AUTO_LOCK                  {0}
ST_FAST_MTEOM                 {0}
ST_SCSI2LOGICAL               {0}
ST_SYSV                       {0}
```

The following parameters are defined in linux/drivers/scsi/st.c

```
ST_TIMEOUT                    {900*HZ}
ST_LONG_TIMEOUT               {14000*HZ}
```

## st ioctls and user interface

The Linux tape interface is defined in `/usr/src/linux/include/linux/mtio.h` .

The following ioctl()s are listed in alphabetical order with a brief explanation to the right. [See st documentation (especially **man 4 st**) for more details.]

```
MTIOCTOP   [execute tape commands and set drive/driver options]
MTIOCGET   [get the status of the drive]
MTIOCPOS   [get the current tape location]

open()   O_RDONLY, O_RDWR
close()
read()
write()
ioctl()
```

# sg driver

The following header files in the kernel source are relevant to the sg driver:

```
/usr/src/linux/include/scsi/sg.h
```

As pointed out in the section called "Mid level" this is best included in applications by using:

```
#include <linux/../scsi/sg.h>
```

# sg compile options

Here are some defines from the sg.h file that the user could conceivably want to change. The current default values are shown in braces on the right:

```
SG_SCATTER_SZ            {32768}
SG_DEF_RESERVED_SIZE     {SG_SCATTER_SZ}
SG_DEF_FORCE_LOW_DMA     {0}
SG_DEF_FORCE_PACK_ID     {0}
SG_DEF_KEP_ORPHAN        {0}
SG_MAX_QUEUE             {16}
SG_DEFAULT_RETRIES       {1}  # i.e. don't retry
SG_BIG_BUFF              {SG_DEF_RESERVED_SIZE}
SG_DEFAULT_TIMEOUT       {60 seconds}
SG_DEF_COMMAND_Q         {0 *}
SG_DEF_UNDERRUN_FLAG     {0}

* The per file descriptor copy of this flips to 1 (thus
  allowing command queuing) as soon as a write() based
  on the newer sg_io_hdr structure is detected.
```

# sg ioctls and user interface

The following ioctl()s are listed in alphabetical order with a brief explanation to the right. [See sg documentation for more details.]

```
SG_EMULATED_HOST      [indicate if adapter is ide-scsi]
SG_GET_COMMAND_Q      [command queuing flag state]
SG_GET_KEEP_ORPHAN    [interrupted SG_IO keep orphan flag state]
SG_GET_LOW_DMA        ["low dma flag" (<= 16 MB on i386) state]
SG_GET_NUM_WAITING    [number of responses waiting to be read()]
SG_GET_PACK_ID        [pack_id of next to read() response
                       (-1 if none)]
SG_GET_REQUEST_TABLE  [yields array of requests being processed]
SG_GET_RESERVED_SIZE  [current size of reserved buffer]
SG_GET_SCSI_ID        [a little more info than the mid level's
                       SCSI_IOCTL_GET_IDLUN ioctl]
SG_GET_SG_TABLESIZE   [max entries in host's scatter gather table]
SG_GET_TIMEOUT        [yields timeout (unit: jiffies
                       (10ms on i386))]
SG_GET_TRANSFORM      [state of ide-scsi's transform flag]
SG_IO                 [send given SCSI command and wait for
                       response]
SG_NEXT_CMD_LEN       [change command length of next command]
SG_SCSI_RESET         [send a SCSI bus, device or host reset]
SG_SET_COMMAND_Q      [set command queuing state {old=0, new=1}]
SG_SET_DEBUG          [set debug level {0}]
SG_SET_KEEP_ORPHAN    [set SG_IO's keep orphan flag {0}]
SG_SET_FORCE_LOW_DMA  [force DMA buffer low (<= 16 MB on i386)
                       {0}]
SG_SET_FORCE_PACK_ID  [so read() can fetch by pack_id {0}]
SG_SET_RESERVED_SIZE  [change default buffer size
                       {SG_DEF_RESERVED_SIZE}]
SG_SET_TIMEOUT        [change current timeout {60 secs} ]
SG_SET_TRANSFORM      [set ide-scsi's ATAPI transform flag {0}]

open()    [recognized oflags: O_RDONLY, O_RDWR, O_EXCL,
           O_NONBLOCK]
close()
read()
write()
ioctl()
poll()    [used when in O_NONBLOCK mode]
fasync()  [enables generation of SIGIO signal for read()]
```

# Appendix F. References, Credits and Corrections

**WEB.** The following references are found on the web. Please alert the author if any of these links become stale.

[W1] SCSI (draft) standards, resources: `www.t10.org` [http://www.t10.org]

[W2] Eric Youngdale is the chief architect of the Linux SCSI subsystem: `www.andante.org/scsi.html` [http://www.andante.org/scsi.html]

[W4] The author's scsi generic (sg) site is: `www.torque.net/sg` [http://www.torque.net/sg]. The Linux Documentation Project's site includes the `www.tldp.org/HOWTO/SCSI-Generic-HOW-TO/` [http://www.tldp.org/HOWTO/SCSI-Generic-HOWTO/] . A (possibly later) version of that document can be found at `www.torque.net/sg/p/sg_v3_ho` [http://www.torque.net/sg/p/sg_v3_ho] . The sg_utils and sg3_utils packages, as tarballs and as binary and source rpms can also be found on this page. These packages and others available for the sg driver are discussed at `www.torque.net/sg/u_index.html` [http://www.torque.net/sg/u_index.html].

[W5] Richard Gooch's devfs site: `www.atnf.csiro.au/~rgooch/linux/docs/devfs.html` [http://www.atnf.csiro.au/~rgooch/linux/docs/devfs.html]

[W6] Kurt Garloff's site (including the **scsidev** and the **scsiinfo** utilities): `www.garloff.de/kurt/linux/` [http://www.garloff.de/kurt/linux/]. Kurt also has the damaged media rescue program **dd_rescue** at this site: `www.garloff.de/kurt/linux/ddrescue` [http://www.garloff.de/kurt/linux/ddrescue]

[W7] Drew Eckhardt's SCSI-HOWTO from 1996 (in ASCII): `metalab.unc.edu/pub/Linux/docs/HOWTO/unmaintained/SCSI-HOWTO` [http://metalab.unc.edu/pub/Linux/docs/HOWTO/unmaintained/SCSI-HOWTO]

[W8] Linux Documentation Project (LDP): `tldp.org` [http://tldp.org]

[W9] SCSI Trade Association site has a lot of useful information: `www.scsita.org` [http://www.scsita.org]

[W10] SCSI FAQ site - useful source of information and links: `www.scsifaq.org` [http://www.scsifaq.org]

**NEWSGROUPS.** The following entries are actually reflectors rather than newsgroups. Various web locations archive their contents (e.g. `marc.theaimsgroup.com` [http://marc.theaimsgroup.com]).

[N1] Linux SCSI reflector: < `linux-scsi@vger.kernel.org` >. This is a relatively low volume (circa 200 postings per month) Linux SCSI specific group that many of the SCSI subsystem maintainers monitor.

[N2] Linux kernel reflector: < `linux-kernel@vger.kernel.org` >. This is a relatively high volume (circa 5000 postings per month) group for all aspects of the Linux kernel. The Linux SCSI reflector should be tried first.

**BOOKS.** Here are some books that the author found useful.

[B1] "Linux Device Drivers" Second edition by Alessandro Rubini and Jonathan Corbet [O'Reilly 2001 ISBN 0-596-00008-1] This is a solid text on Linux device drivers including some information on the SCSI

subsystem. It covers the block subsystem well and has many char device driver examples. It has been updated for the Linux 2.4 series kernels and also includes information on the Linux 2.2 and 2.0 series. This book is highly recommended. The authors and the publisher have unselfishly made this book available under the GNU Free Documentation License (version 1.1). It can be found in html at `www.xml.com/ldd/chapter/book` [http://www.xml.com/ldd/chapter/book] .

 [B2] "Running Linux" 3rd edition by M. Welsh, M. K. Dalheimer & L. Kaufman [O'Reilly 1999 ISBN 1-56592-469-X] This is a classic Linux tome which includes some SCSI configuration info.

 [B3] "The Programmer's Guide to SCSI" by Brian Sawert [Addison Wesley 1998 ISBN 0-201-18538-5] This book covers many SCSI topics, including the pass through mechanisms of Linux (sg) and ASPI/ASPI32 as used by Windows.

**CREDITS.**     The author is grateful for the following contributions:

- Kai Mäkisara (st) `<Kai.Makisara at metla dot fi>`

- Jens Axboe (sr) `<axboe at suse dot de>`

- Richard Gooch (devfs) `<rgooch at atnf dot csiro dot au>`

- Tim Waugh (ppa, imm, ppscsi + docbook) `<twaugh at redhat dot com>`

- Gadi Oxman (ide-scsi) `<gadio at netvision dot net dot il>`

**CORRECTIONS and SUGGESTIONS.**     Please send any corrections or suggestions to the author at `<dgilbert at interlog dot com>` or `<dougg at torque dot net>` .