

Large Disk HOWTO

Table of Contents

Large Disk HOWTO	1
<u>Andries Brouwer, aeb@cwi.nl</u>	1
<u>1. Large disks</u>	1
<u>2. Units</u>	2
<u>3. Disk Access</u>	2
<u>3.1 Cylinders, heads and sectors</u>	3
<u>3.2 Sectorsize</u>	3
<u>3.3 Disksize</u>	3
<u>3.4 The 1024 cylinder and 8.5 GB limits</u>	3
<u>3.5 The 137 GB limit</u>	4
<u>4. History of BIOS and IDE limits</u>	4
<u>5. Booting</u>	5
<u>5.1 LILO and the `lba32' and `linear' options</u>	6
<u>5.2 A LILO bug</u>	6
<u>5.3 1024 cylinders is not 1024 cylinders</u>	6
<u>5.4 No 1024 cylinder limit on old machines with IDE</u>	7
<u>5.5 Other boot loaders</u>	7
<u>6. Disk geometry, partitions and `overlap'</u>	7
<u>6.1 The last cylinder</u>	8
<u>6.2 Cylinder boundaries</u>	8
<u>7. Translation and Disk Managers</u>	9
<u>8. Kernel disk translation for IDE disks</u>	10
<u>8.1 EZD</u>	10
<u>8.2 DM6:DDO</u>	10
<u>8.3 DM6:AUX</u>	10
<u>8.4 DM6:MBR</u>	10
<u>8.5 PTBL</u>	10
<u>8.6 Getting rid of a disk manager</u>	11
<u>8.7 Since 2.5.70: boot parameters</u>	11
<u>9. Consequences</u>	11
<u>9.1 Computing LILO parameters</u>	13
<u>10. Details</u>	13
<u>10.1 IDE details - the seven geometries</u>	13
<u> The IDENTIFY DRIVE command</u>	14
<u>10.2 SCSI details</u>	15
<u>11. Clipped disks</u>	17
<u>11.1 The Linux IDE 8 GiB limit</u>	17
<u>11.2 BIOS complications</u>	17
<u>11.3 Jumpers that select the number of heads</u>	18
<u>11.4 Jumpers that clip total capacity</u>	18
<u> Clip to 2.1 GB</u>	18
<u> Clip to 33 GB</u>	18
<u> Maxtor</u>	19
<u> IBM</u>	19
<u> Seagate</u>	19
<u> Maxtor D540X-4K</u>	20
<u> Western Digital</u>	20
<u>11.5 READ NATIVE MAX ADDRESS / SET MAX ADDRESS</u>	20

Table of Contents

Large Disk HOWTO

<u>11.6 CONFIG IDEDISK STROKE</u>	21
<u>12. The Linux 65535 cylinder limit</u>	21
<u>12.1 IDE problems with 34+ GB disks</u>	21
<u>13. Extended and logical partitions</u>	22
<u>14. Problem solving</u>	23
<u>14.1 Problem: My IDE disk gets a bad geometry when I boot from SCSI</u>	23
<u>14.2 Nonproblem: Identical disks have different geometry?</u>	23
<u>14.3 Problem: 2.4 and 2.6 report different geometries? 2.6 reports the wrong geometry? 2.6 reports no geometry at all?</u>	24
<u>14.4 Nonproblem: fdisk sees much more room than df?</u>	24

Large Disk HOWTO

Andries Brouwer, aeb@cwi.nl

v2.5, 2004-11-01

All about disk geometry and the 1024 cylinder and other limits for disks.

For the most recent version of this text, see www.win.tue.nl.

1. Large disks

You got a new disk. What to do? Well, on the software side: use `fdisk` or `cdisk` to create partitions, and then `mke2fs` or `mkreiserfs` or so to create a filesystem, and then `mount` to attach the new filesystem to the big file hierarchy. Make sure you have relatively recent versions of these utilities - often old versions have problems handling large disks.

You need not read this HOWTO since there are *no* problems with large hard disks these days.

Long ago, disks were large when they had a capacity larger than 528 MB, or than 8.4 GB, or than 33.8 GB. These days the interesting limit is 137 GB. In all cases, sufficiently recent Linux kernels handle the disk fine.

Sometimes booting requires some care, since Linux cannot help you when it isn't running yet. But again, with a sufficiently recent BIOS and boot loader there are no problems. Most of the text below will treat the cases of (i) ancient hardware, (ii) broken hardware or BIOS, (iii) several operating systems on the same disk, (iv) booting old systems.

Advice

For large SCSI disks: Linux has supported them from very early on. No action required.

For large IDE disks (over 8.4 GB): make sure your kernel is 2.0.34 or later.

For large IDE disks (over 33.8 GB): make sure your kernel is 2.0.39/2.2.14/2.3.21 or later.

For large IDE disks (over 137 GB): make sure your kernel is 2.4.19/2.5.3 or later.

If the kernel boots fine, and the boot messages indicate that it recognizes the disk correctly, but there are problems with utilities, upgrade the utilities.

If LILO hangs at boot time, make sure you have version 21.4 or later, and specify the keyword `lba32` in the configuration file `/etc/lilo.conf`. With an older version of LILO, try both with and without the `linear` keyword.

There may be geometry problems that can be solved by giving an explicit geometry to kernel/LILO/fdisk.

If you have an old `fdisk` and it warns about overlapping partitions: ignore the warnings, or check using `cdisk` that really all is well.

Large Disk HOWTO

For HPT366, see the [Linux HPT366 HOWTO](#).

If at boot time the kernel cannot read the partition table, consider the possibility that UDMA66 was selected while the controller or the cable or the disk drive did not support UDMA66. In such a case every attempt to read will fail, and reading the partition table is the first thing the kernel does. Make sure no UDMA66 is used.

If the BIOS hangs at boot time because of a large disk, and flashing a newer version is not an option, take the disk out of the BIOS setup. If you have to boot from the disk, look whether a capacity clipping jumper helps.

If you think something is wrong with the size of your disk, make sure that you are not confusing binary and decimal [units](#), and realize that the free space that `df` reports on an empty disk is a few percent smaller than the partition size, because there is administrative overhead. Software that does not understand 48-bit addressing will view a 137+ GB disk as having a capacity of 137 GB. When a capacity clipping [jumper](#) is present, a larger disk may have been clipped to 33 GB or to 2 GB.

If for a removable drive the kernel reports two different sizes, then one is found from the drive, and the other from the disk/floppy. This second value will be zero when the drive has no media.

Now, if you still think there are problems, or just are curious, read on.

Below a rather detailed description of all relevant details. I used kernel version 2.0.8 source as a reference. Other versions may differ a bit.

2. Units

A kilobyte (kB) is 1000 bytes. A megabyte (MB) is 1000 kB. A gigabyte (GB) is 1000 MB. A terabyte (TB) is 1000 GB. This is the [SI norm](#). However, there are people that use 1 MB=1024000 bytes and talk about 1.44 MB floppies, and people who think that 1 MB=1048576 bytes. Here I follow the [recent standard](#) and write Ki, Mi, Gi, Ti for the binary units, so that these floppies are 1440 KiB (1.47 MB, 1.41 MiB), 1 MiB is 1048576 bytes (1.05 MB), 1 GiB is 1073741824 bytes (1.07 GB) and 1 TiB is 1099511627776 bytes (1.1 TB).

Quite correctly, the disk drive manufacturers follow the SI norm and use the decimal units. However, Linux kernel boot messages (for not-so-recent kernels) and some old fdisk-type programs use the symbols MB and GB for binary, or mixed binary-decimal units. So, before you think your disk is smaller than was promised when you bought it, compute first the actual size in decimal units (or just in bytes).

Concerning terminology and abbreviation for binary units, [Knuth](#) has an alternative [proposal](#), namely to use KKB, MMB, GGB, TTB, PPB, EEB, ZZB, YYB and to call these *large kilobyte*, *large megabyte*, ... *large yottabyte*. He writes: 'Notice that doubling the letter connotes both binary-ness and large-ness.' This is a good proposal - 'large gigabyte' sounds better than 'gibibyte'. For our purposes however the only important thing is to stress that a megabyte has precisely 1000000 bytes, and that some other term and abbreviation is required if you mean something else.

3. Disk Access

Disk access is done in units called *sectors*. In order to read or write something from or to the disk, we have to specify the position on the disk, for example by giving the sector number. If the disk is a SCSI disk, then this sector number goes directly into the SCSI command and is understood by the disk. If the disk is an IDE disk using LBA, then precisely the same holds. But if the disk is old, RLL or MFM or IDE from before the LBA times, then the disk hardware expects a triple (cylinder,head,sector) to designate the desired spot on the disk.

3.1 Cylinders, heads and sectors

A disk has sectors numbered 0, 1, 2, ... This is called *LBA addressing*.

In ancient times, before the advent of IDE disks, disks had a *geometry* described by three constants C, H, S: the number of cylinders, the number of heads, the number of sectors per track. The address of a sector was given by three numbers: *c, h, s*: the cylinder number (between 0 and C-1), the head number (between 0 and H-1), and the sector number within the track (between 1 and S), where for some mysterious reason *c* and *h* count from 0, but *s* counts from 1. This is called *CHS addressing*.

No disk manufactured less than ten years ago has a geometry, but this ancient 3D sector addressing is still used by the INT13 BIOS interface (with fantasy numbers C, H, S unrelated to any physical reality).

The correspondence between the linear numbering and this 3D notation is as follows: for a disk with C cylinders, H heads and S sectors/track position (*c,h,s*) in 3D or CHS notation is the same as position $c * H * S + h * S + (s-1)$ in linear or LBA notation.

Consequently, in order to access a very old non-SCSI disk, we need to know its *geometry*, that is, the values of C, H and S. (And if you don't know, there is a lot of good information on www.thetechpage.com.)

3.2 Sectorsize

In the present text a sector has 512 bytes. This is almost always true, but for example certain MO disks use a sectorsize of 2048 bytes, and all capacities given below must be multiplied by four. (When using `fdisk` on such disks, make sure you have version 2.9i or later, and give the ``-b 2048'` option.)

3.3 Disksize

A disk with C cylinders, H heads and S sectors per track has $C * H * S$ sectors in all, and can store $C * H * S * 512$ bytes. For example, if the disk label says C/H/S=4092/16/63 then the disk has $4092 * 16 * 63 = 4124736$ sectors, and can hold $4124736 * 512 = 2111864832$ bytes (2.11 GB). There is an industry convention to give C/H/S=16383/16/63 for disks larger than 8.4 GB, and the disk size can no longer be read off from the C/H/S values reported by the disk.

3.4 The 1024 cylinder and 8.5 GB limits

The old INT13 BIOS interface to disk I/O uses 24 bits to address a sector: 10 bits for the cylinder, 8 bits for the head, and 6 bits for the sector number within the track (counting from 1). This means that this interface cannot address more than $1024 * 256 * 63$ sectors, which is 8.5 GB (with 512-byte sectors). And if the (fantasy) geometry specified for the disk has fewer than 1024 cylinders, or 256 heads, or 63 sectors per track, then this limit will be less.

(More precisely: with INT 13, AH selects the function to perform, CH is the low 8 bits of the cylinder number, CL has in bits 7-6 the high two bits of the cylinder number and in bits 5-0 the sector number, DH is the head number, and DL is the drive number (80h or 81h). This explains part of the layout of the partition table.)

This state of affairs was rectified when the so-called Extended INT13 functions were introduced. A modern BIOS has no problems accessing large disks.

Large Disk HOWTO

(More precisely: DS:SI points at a 16-byte Disk Address Packet that contains an 8-byte starting absolute block number.)

Linux does not use the BIOS, so does (and did) not have this problem.

However, this geometry stuff plays a role in the interpretation of partition tables, so if Linux shares a disk with for example DOS, then it needs to know what geometry DOS will think the disk has. It also plays a role at boot time, where the BIOS has to load a boot loader, and the boot loader has to load the operating system.

3.5 The 137 GB limit

The old ATA standard describes how to address a sector on an IDE disk using 28 bits (8 bits for the sector, 4 for the head, 16 for the cylinder). This means that an IDE disk can have at most 2^{28} addressable sectors. With 512-byte sectors this is 2^{37} bytes, that is, 137.4 GB.

The ATA-6 standard includes a specification how to address past this 2^{28} sector boundary. The new standard allows addressing of 2^{48} sectors. There is support in recent Linux kernels that have incorporated Andre Hedrick's IDE patch, for example 2.4.18-pre7-ac3 and 2.5.3.

Maxtor sells 160 GB IDE disks since Fall 2001. An old kernel will treat such disks as 137.4 GB disks.

4. History of BIOS and IDE limits

ATA Specification (for IDE disks) - the 137 GB limit

At most 65536 cylinders (numbered 0-65535), 16 heads (numbered 0-15), 255 sectors/track (numbered 1-255), for a maximum total capacity of 267386880 sectors (of 512 bytes each), that is, 136902082560 bytes (137 GB). In Sept 2001, the first drives larger than this (160 GB Maxtor Diamondmax) appeared.

BIOS Int 13 - the 8.5 GB limit

At most 1024 cylinders (numbered 0-1023), 256 heads (numbered 0-255), 63 sectors/track (numbered 1-63) for a maximum total capacity of 8455716864 bytes (8.5 GB). This is a serious limitation today. It means that DOS cannot use present day large disks.

The 528 MB limit

If the same values for c,h,s are used for the BIOS Int 13 call and for the IDE disk I/O, then both limitations combine, and one can use at most 1024 cylinders, 16 heads, 63 sectors/track, for a maximum total capacity of 528482304 bytes (528MB), the infamous 504 MiB limit for DOS with an old BIOS. This started being a problem around 1993, and people resorted to all kinds of trickery, both in hardware (LBA), in firmware (translating BIOS), and in software (disk managers). The concept of 'translation' was invented (1994): a BIOS could use one geometry while talking to the drive, and another, fake, geometry while talking to DOS, and translate between the two.

The 2.1 GB limit (April 1996)

Some older BIOSes only allocate 12 bits for the field in CMOS RAM that gives the number of cylinders. Consequently, this number can be at most 4095, and only $4095 * 16 * 63 * 512 = 2113413120$ bytes are accessible. The effect of having a larger disk would be a hang at boot time. This made disks with geometry 4092/16/63 rather popular. And still today many large disk drives come with a jumper to make them appear 4092/16/63. See also [over2gb.htm](#). [Other BIOSes](#) would not hang but just detect a much smaller disk, like 429 MB instead of 2.5 GB.

The 3.2 GB limit

There was a bug in the Phoenix 4.03 and 4.04 BIOS firmware that would cause the system to lock up in the CMOS setup for drives with a capacity over 3277 MB. See [over3gb.htm](#).

Large Disk HOWTO

The 4.2 GB limit (Feb 1997)

Simple BIOS translation (ECHS=Extended CHS, sometimes called 'Large disk support' or just 'Large') works by repeatedly doubling the number of heads and halving the number of cylinders shown to DOS, until the number of cylinders is at most 1024. Now DOS and Windows 95 cannot handle 256 heads, and in the common case that the disk reports 16 heads, this means that this simple mechanism only works up to $8192 * 16 * 63 * 512 = 4227858432$ bytes (with a fake geometry with 1024 cylinders, 128 heads, 63 sectors/track). Note that ECHS does not change the number of sectors per track, so if that is not 63, the limit will be lower. See [over4gb.htm](#).

The 7.9 GB limit

Slightly smarter BIOSes avoid the previous problem by first adjusting the number of heads to 15 ('revised ECHS'), so that a fake geometry with 240 heads can be obtained, good for $1024 * 240 * 63 * 512 = 7927234560$ bytes.

The 8.4 GB limit

Finally, if the BIOS does all it can to make this translation a success, and uses 255 heads and 63 sectors/track ('assisted LBA' or just 'LBA') it may reach $1024 * 255 * 63 * 512 = 8422686720$ bytes, slightly less than the earlier 8.5 GB limit because the geometries with 256 heads must be avoided. (This translation will use for the number of heads the first value H in the sequence 16, 32, 64, 128, 255 for which the total disk capacity fits in $1024 * H * 63 * 512$, and then computes the number of cylinders C as total capacity divided by $(H * 63 * 512)$.)

The 33.8 GB limit (August 1999)

The next hurdle comes with a size over 33.8 GB. The problem is that with the default 16 heads and 63 sectors/track this corresponds to a number of cylinders of more than 65535, which does not fit into a short. Many BIOSes couldn't handle such disks. (See, e.g., [Asus upgrades](#) for new flash images that work.) Linux kernels older than 2.2.14 / 2.3.21 need a patch. See [IDE problems with 34+ GB disks](#) below.

The 137 GB limit (Sept 2001)

As mentioned above, the old ATA protocol uses $16+4+8 = 28$ bits to specify the sector number, and hence cannot address more than 2^{28} sectors. ATA-6 describes an extension that allows the addressing of 2^{48} sectors, a million times as much. There is support in very recent kernels.

The 2 TiB limit

With 32-bit sector numbers, one can address 2 TiB. A lot of software will have to be rewritten once disks get larger.

Hard drives over 8.4 GB are supposed to report their geometry as 16383/16/63. This in effect means that the 'geometry' is obsolete, and the total disk size can no longer be computed from the geometry, but is found in the LBA capacity field returned by the [IDENTIFY command](#). Hard drives over 137.4 GB are supposed to report an LBA capacity of $0xffffffff = 268435455$ sectors (137438952960 bytes). Now the actual disk size is found in the new 48-capacity field.

5. Booting

When the system is booted, the BIOS reads sector 0 (known as the MBR - the Master Boot Record) from the first disk (or from floppy or CDROM), and jumps to the code found there - usually some bootstrap loader. These small bootstrap programs found there typically have no own disk drivers and use BIOS services. This means that a Linux kernel can only be booted when it is entirely located within the first 1024 cylinders, unless you both have a modern BIOS (a BIOS that supports the Extended INT13 functions), and a modern bootloader (a bootloader that uses these functions when available).

This problem (if it is a problem) is very easily solved: make sure that the kernel (and perhaps other files used during bootup, such as LILO map files) are located on a partition that is entirely contained in the first 1024

Large Disk HOWTO

cylinders of a disk that the BIOS can access - probably this means the first or second disk.

Thus: create a small partition, say 10 MB large, so that there is room for a handful of kernels, making sure that it is entirely contained within the first 1024 cylinders of the first or second disk. Mount it on `/boot` so that LILO will put its stuff there.

Most systems from 1998 or later will have a modern BIOS.

5.1 LILO and the ``lba32'` and ``linear'` options

Executive summary: If you use LILO as boot loader, make sure you have LILO version 21.4 or later. (It can be found at <ftp://metalab.unc.edu/pub/Linux/system/boot/lilo/>.) Always use the `lba32` option.

An invocation of `/sbin/lilo` (the boot map installer) stores a list of addresses in the boot map, so that LILO (the boot loader) knows from where to read the kernel image. By default these addresses are stored in (c,h,s) form, and ordinary INT13 calls are used at boot time.

When the configuration file specifies `lba32` or `linear`, linear addresses are stored. With `lba32` also linear addresses are used at boot time, when the BIOS supports extended INT13. With `linear`, or with an old BIOS, these linear addresses are converted back to (c,h,s) form, and ordinary INT13 calls are used.

Thus, with `lba32` there are no geometry problems and there is no 1024 cylinder limit. Without it there is a 1024 cylinder limit. What about the geometry?

The boot loader and the BIOS must agree as to the disk geometry. `/sbin/lilo` asks the kernel for the geometry, but there is no guarantee that the Linux kernel geometry coincides with what the BIOS will use. Thus, often the geometry supplied by the kernel is worthless. In such cases it helps to give LILO the ``linear'` option. The advantage is that the Linux kernel idea of the geometry no longer plays a role. The disadvantage is that `lilo` cannot warn you when part of the kernel was stored above the 1024 cylinder limit, and you may end up with a system that does not boot.

5.2 A LILO bug

With LILO versions below v21 there is another disadvantage: the address conversion done at boot time has a bug: when `c*H` is 65536 or more, overflow occurs in the computation. For `H` larger than 64 this causes a stricter limit on `c` than the well-known `c < 1024`; for example, with `H=255` and an old LILO one must have `c < 258`. (`c`=cylinder where kernel image lives, `H`=number of heads of disk)

5.3 1024 cylinders is not 1024 cylinders

Tim Williams writes: 'I had my Linux partition within the first 1024 cylinders and still it wouldnt boot. First when I moved it below 1 GB did things work.' How can that be? Well, this was a SCSI disk with AHA2940UW controller which uses either `H=64`, `S=32` (that is, cylinders of 1 MiB = 1.05 MB), or `H=255`, `S=63` (that is, cylinders of 8.2 MB), depending on setup options in firmware and BIOS. No doubt the BIOS assumed the former, so that the 1024 cylinder limit was found at 1 GiB, while Linux used the latter and LILO thought that this limit was at 8.4 GB.

5.4 No 1024 cylinder limit on old machines with IDE

The nuni boot loader does not use BIOS services but accesses IDE drives directly. So, one can put it on a floppy or in the MBR and boot from anywhere on any IDE drive (not only from the first two). Find it at [//metalab.unc.edu/pub/Linux/system/boot/loaders/](http://metalab.unc.edu/pub/Linux/system/boot/loaders/).

5.5 Other boot loaders

LILO is a bit fragile, it requires the discipline of running `/sbin/lilo` each time one installs a new kernel. Some other boot loaders do not have this disadvantage. Especially `grub` is popular these days; a major disadvantage is that it does not support the `lilo -R label` function.

6. Disk geometry, partitions and `overlap'

If you have several operating systems on your disks, then each uses one or more disk partitions. A disagreement on where these partitions are may have catastrophic consequences.

The MBR contains a *partition table* describing where the (primary) partitions are. There are 4 table entries, for 4 primary partitions, and each looks like

```
struct partition {
    char active;      /* 0x80: bootable, 0: not bootable */
    char begin[3];   /* CHS for first sector */
    char type;
    char end[3];     /* CHS for last sector */
    int start;       /* 32 bit sector number (counting from 0) */
    int length;      /* 32 bit number of sectors */
};
```

(where CHS stands for Cylinder/Head/Sector).

This information is redundant: the location of a partition is given both by the 24-bit `begin` and `end` fields, and by the 32-bit `start` and `length` fields.

Linux only uses the `start` and `length` fields, and can therefore handle partitions of not more than 2^{32} sectors, that is, partitions of at most 2 TiB. That is twelve times larger than the disks available today, so maybe it will be enough for the next five years or so. (So, partitions can be very large, but there is a serious restriction in that a file in an ext2 filesystem on hardware with 32-bit integers cannot be larger than 2 GiB.)

DOS uses the `begin` and `end` fields, and uses the BIOS INT13 call to access the disk, and can therefore only handle disks of not more than 8.4 GB, even with a translating BIOS. (Partitions cannot be larger than 2.1 GB because of restrictions of the FAT16 file system.) The same holds for Windows 3.11 and WfWG and Windows NT 3.*.

Windows 95 has support for the Extended INT13 interface, and uses special partition types (c, e, f instead of b, 6, 5) to indicate that a partition should be accessed in this way. When these partition types are used, the `begin` and `end` fields contain dummy information (1023/255/63). Windows 95 OSR2 introduces the FAT32 file system (partition type b or c), that allows partitions of size at most 2 TiB.

What is this nonsense you get from `fdisk` about `overlapping' partitions, when in fact nothing is wrong? Well - there is something `wrong': if you look at the `begin` and `end` fields of such partitions, as DOS does,

Large Disk HOWTO

they overlap. (And that cannot be corrected, because these fields cannot store cylinder numbers above 1024 - there will always be `overlap' as soon as you have more than 1024 cylinders.) However, if you look at the `start` and `length` fields, as Linux does, and as Windows 95 does in the case of partitions with partition type `c`, `e` or `f`, then all is well. So, ignore these warnings when `cfdisk` is satisfied and you have a Linux-only disk. Be careful when the disk is shared with DOS. Use the commands `cfdisk -Ps /dev/hdx` and `cfdisk -Pt /dev/hdx` to look at the partition table of `/dev/hdx`.

6.1 The last cylinder

Many old IBM PS/2 systems used disks with a defect map written to the end of the disk. (Bit 0x20 in the control word of the [disk parameter table](#) is set.) Therefore, FDISK would not use the last cylinder. Just to be sure, the BIOS often already reports the size of the disk as one cylinder smaller than reality, and that may mean that two cylinders are lost. Newer BIOSes have several disk size reporting functions, where internally one calls the other. When both subtract 1 for this reserved cylinder and also FDISK does so, then one may lose three cylinders. These days all of this is irrelevant, but this may provide an explanation if one observes that different utilities have slightly different opinions about the disk size.

6.2 Cylinder boundaries

A well-known claim says that partitions should start and end at cylinder boundaries.

Since "disk geometry" is something without objective existence, different operating systems will invent different geometries for the same disk. One often sees a translated geometry like `*/255/63` used by one and an untranslated geometry like `*/16/63` used by another OS. (People tell me Windows NT uses `*/64/32` while Windows 2K uses `*/255/63`.) Thus, it may be impossible to align partitions to cylinder boundaries according to each of the various ideas about the size of a cylinder that one's systems have. Also different Linux kernels may assign different geometries to the same disk. Also, enabling or disabling the BIOS of a SCSI card may change the fake geometry of the connected SCSI disks.

Fortunately, for Linux there is no alignment requirement at all. (Except that some semi-broken installation software likes to be very sure that all is OK; thus, it may be impossible to install RedHat 7.1 on a disk with unaligned partitions because DiskDruid is unhappy.)

People report that it is easy to create nonaligned partitions in Windows NT, without any noticeable bad effects.

But MSDOS 6.22 has an alignment requirement. Extended partition sectors that are not on a cylinder boundary are ignored by its FDISK. The system itself is happy with any alignment, but interprets relative starting addresses as if relative to an aligned address: The starting address of a logical partition is given relative not to the address of the extended partition sector that describes it, but relative to the start of the cylinder that contains that sector. (So, it is not surprising that also PartitionMagic requires alignment.)

What is the definition of alignment? MSDOS 6.22 FDISK will do the following: 1. If the first sector of the cylinder is a partition table sector, then the rest of the track is unused, and the partition starts with the next track. This applies to sector 0 (the MBR) and the partition table sectors preceding logical partitions. 2. Otherwise, the partition starts at the first sector of the cylinder. Also the extended partition starts at a cylinder boundary. The `cfdisk` man page says that old versions of DOS did not align partitions.

Use of partition type 85 for the extended partition makes it invisible to DOS, making sure that only Linux will look inside.

Large Disk HOWTO

As an aside: on a Sparc, the boot partition must start on a cylinder boundary (but there is no requirement on the end).

7. Translation and Disk Managers

Disk geometry (with heads, cylinders and tracks) is something from the age of MFM and RLL. In those days it corresponded to a physical reality. Nowadays, with IDE or SCSI, nobody is interested in what the 'real' geometry of a disk is. Indeed, the number of sectors per track is variable - there are more sectors per track close to the outer rim of the disk - so there is no 'real' number of sectors per track. Quite the contrary: the IDE command INITIALIZE DRIVE PARAMETERS (91h) serves to tell the disk how many heads and sectors per track it is supposed to have today. It is quite normal to see a large modern disk that has 2 heads report 15 or 16 heads to the BIOS, while the BIOS may again report 255 heads to user software.

For the user it is best to regard a disk as just a linear array of sectors numbered 0, 1, ..., and leave it to the firmware to find out where a given sector lives on the disk. This linear numbering is called LBA.

So now the conceptual picture is the following. DOS, or some boot loader, talks to the BIOS, using (c,h,s) notation. The BIOS converts (c,h,s) to LBA notation using the fake geometry that the user is using. If the disk accepts LBA then this value is used for disk I/O. Otherwise, it is converted back to (c',h',s') using the geometry that the disk uses today, and that is used for disk I/O.

Note that there is a bit of confusion in the use of the expression 'LBA': As a term describing disk capabilities it means 'Linear Block Addressing' (as opposed to CHS Addressing). As a term in the BIOS Setup, it describes a translation scheme sometimes called 'assisted LBA' - see above under '[The 8.4 GB limit](#)'.

Something similar works when the firmware doesn't speak LBA but the BIOS knows about translation. (In the setup this is often indicated as 'Large'.) Now the BIOS will present a geometry (C,H,S) to the operating system, and use (C',H',S') while talking to the disk controller. Usually $S = S'$, $C = C'/N$ and $H = H' * N$, where N is the smallest power of two that will ensure $C' \leq 1024$ (so that least capacity is wasted by the rounding down in $C' = C/N$). Again, this allows access of up to 8.4 GB (7.8 GiB).

(The third setup option usually is 'Normal', where no translation is involved.)

If a BIOS does not know about 'Large' or 'LBA', then there are software solutions around. Disk Managers like OnTrack or EZ-Drive replace the BIOS disk handling routines by their own. Often this is accomplished by having the disk manager code live in the MBR and subsequent sectors (OnTrack calls this code DDO: Dynamic Drive Overlay), so that it is booted before any other operating system. That is why one may have problems when booting from a floppy when a Disk Manager has been installed.

The effect is more or less the same as with a translating BIOS - but especially when running several different operating systems on the same disk, disk managers can cause a lot of trouble.

Linux did support OnTrack Disk Manager since version 1.3.14, and EZ-Drive since version 1.3.29. Some more details are given in the next section.

In 2.5.70 the automatic disk manager support was removed. Instead, two boot options were added: "hda=remap" to do the EZ-Drive remapping of sector 0 to sector 1, and "hda=remap63" to do the OnTrack Disk Manager shift over 63 sectors.

8. Kernel disk translation for IDE disks

If the Linux kernel detects the presence of some disk manager on an IDE disk, it will try to remap the disk in the same way this disk manager would have done, so that Linux sees the same disk partitioning as for example DOS with OnTrack or EZ-Drive. However, NO remapping is done when a geometry was specified on the command line - so a ``hd=cyls, heads, secs'` command line option might well kill compatibility with a disk manager.

If you are hit by this, and know someone who can compile a new kernel for you, find the file `linux/drivers/block/ide.c` and remove in the routine `ide_xlate_1024()` the test `if (drive->forced_geom) { ...; return 0; }`.

The remapping is done by trying 4, 8, 16, 32, 64, 128, 255 heads (keeping H*C constant) until either $C \leq 1024$ or $H = 255$.

The details are as follows - subsection headers are the strings appearing in the corresponding boot messages. Here and everywhere else in this text partition types are given in hexadecimal.

8.1 EZD

EZ-Drive is detected by the fact that the first primary partition has type 55. The geometry is remapped as described above, and the partition table from sector 0 is discarded - instead the partition table is read from sector 1. Disk block numbers are not changed, but writes to sector 0 are redirected to sector 1. This behaviour can be changed by recompiling the kernel with `#define FAKE_FDISK_FOR_EZDRIVE 0` in `ide.c`.

8.2 DM6:DDO

OnTrack DiskManager (on the first disk) is detected by the fact that the first primary partition has type 54. The geometry is remapped as described above and the entire disk is shifted by 63 sectors (so that the old sector 63 becomes sector 0). Afterwards a new MBR (with partition table) is read from the new sector 0. Of course this shift is to make room for the DDO - that is why there is no shift on other disks.

8.3 DM6:AUX

OnTrack DiskManager (on other disks) is detected by the fact that the first primary partition has type 51 or 53. The geometry is remapped as described above.

8.4 DM6:MBR

An older version of OnTrack DiskManager is detected not by partition type, but by signature. (Test whether the offset found in bytes 2 and 3 of the MBR is not more than 430, and the short found at this offset equals 0x55AA, and is followed by an odd byte.) Again the geometry is remapped as above.

8.5 PTBL

Finally, there is a test that tries to deduce a translation from the `start` and `end` values of the primary partitions: If some partition has start and end sector number 1 and 63, respectively, and end heads 31, 63, 127

or 254, then, since it is customary to end partitions on a cylinder boundary, and since moreover the IDE interface uses at most 16 heads, it is conjectured that a BIOS translation is active, and the geometry is remapped to use 32, 64, 128 or 255 heads, respectively. However, no remapping is done when the current idea of the geometry already has 63 sectors per track and at least as many heads (since this probably means that a remapping was done already).

8.6 Getting rid of a disk manager

When Linux detects OnTrack Disk Manager, it will shift all disk accesses by 63 sectors. Similarly, when Linux detects EZ-Drive, it shifts all accesses of sector 0 to sector 1. This means that it may be difficult to get rid of these disk managers. Most disk managers have an uninstall option, but if you need to remove some disk manager an approach that often works is to give an explicit disk geometry on the command line. Now Linux skips the `ide_xlate_1024()` routine, and one can wipe out the partition table with disk manager (and probably lose access to all disk data) with the command

```
dd if=/dev/zero of=/dev/hdx bs=512 count=1
```

The details depend a little on kernel version. Recent kernels (since 2.3.21) recognize boot parameters like "hda=remap" and "hdb=noremap", so that it is possible to get or avoid the EZD shift regardless of the contents of the partition table. The "hdX=noremap" boot parameter also avoids the OnTrack Disk Manager shift.

8.7 Since 2.5.70: boot parameters

In 2.5.70 the automatic disk manager support was removed. Instead, two boot options were added: "hda=remap" to do the EZ-Drive remapping of sector 0 to sector 1, and "hda=remap63" to do the OnTrack Disk Manager shift over 63 sectors.

This also means that it no longer is a problem to get rid of a disk manager.

9. Consequences

What does all of this mean? For Linux users only one thing: that they must make sure that LILO and `fdisk` use the right geometry where 'right' is defined for `fdisk` as the geometry used by the other operating systems on the same disk, and for LILO as the geometry that will enable successful interaction with the BIOS at boot time. (Usually these two coincide.)

How does `fdisk` know about the geometry? There are three sources of information. First, if the user has specified the geometry interactively or on the command line, then we take the user input. Second, if it is possible to guess the geometry used from the partition table, then we use that. Third, when nothing else is available, `fdisk` asks the kernel, using the `HDIO_GETGEO` ioctl.

How does LILO know about the geometry? It asks the kernel, using the `HDIO_GETGEO` ioctl. But the user can override the geometry using the `'disk='` option in `/etc/lilo.conf` (see `lilo.conf(5)`). One may also give the `linear` option to LILO, and it will store LBA addresses instead of CHS addresses in its map file, and find out of the geometry to use at boot time (by using INT 13 Function 8 to ask for the drive geometry).

How does the kernel know what to answer? Well, first of all, the user may have specified an explicit geometry with a `'hda=cyls,heads,secs'` kernel command line option (see `bootparam(7)`), perhaps by hand, or by asking the boot loader to supply such an option to the kernel. For example, one can tell LILO to supply such an

Large Disk HOWTO

option by adding an `append = "hda=cyls,heads,secs"` line in `/etc/lilo.conf` (see `lilo.conf(5)`). And otherwise the kernel will guess, possibly using values obtained from the BIOS or the hardware.

Note that values guessed by the kernel are very unreliable. The kernel does not have a good way of finding out what values the BIOS uses, or indeed whether the disk is known to the BIOS at all.

It is possible (since Linux 2.1.79) to change the kernel's ideas about the geometry by using the `/proc` filesystem. For example

```
# sfdisk -g /dev/hdc
/dev/hdc: 4441 cylinders, 255 heads, 63 sectors/track
# cd /proc/ide/ide1/hdc
# echo bios_cyl:17418 bios_head:128 bios_sect:32 > settings
# sfdisk -g /dev/hdc
/dev/hdc: 17418 cylinders, 128 heads, 32 sectors/track
#
```

This is especially useful if you need so many boot parameters that you overflow LILO's (very limited) command line length. (It also helps if you want to influence a utility that gets its idea of the geometry from the kernel via the `HDIO_GETGEO` ioctl.)

Since Linux 2.6.5 the kernel will (when compiled with `CONFIG_EDD`) ask the BIOS for `legacy_cylinders`, `legacy_heads`, `legacy_sectors` using `INT 13/AH=08h`. The values obtained are made available in `/sys/firmware/edd/int13_dev{80,81,82,83}/legacy_*`. In 2.6.5 the files were `legacy_{cylinders,heads,sectors}` (with contents in hex, e.g. `0xfe` for 254), but those names are confusing, and in 2.6.7 they were changed to `legacy_max_cylinder`, `legacy_max_head`, `legacy_sectors_per_track` (with contents in decimal). A geometry like `C/H/S=1000/255/63` is found here as `999, 254, 63`.

```
# insmod edd.ko
# cd /sys/firmware/edd/int13_dev83
# cat legacy_max_head
254
# cat sectors
120064896
#
```

Thus, we see here a disk with 255 heads and 120064896 sectors in all. Careful comparison shows that this is `/dev/hdf`.

How does the BIOS know about the geometry? The user may have specified it in the CMOS setup. Or the geometry is read from the disk, and possibly translated as specified in the setup. In the case of SCSI disks, where no geometry exists, the geometry that the BIOS has to invent can often be specified by jumpers or setup options. (For example, Adaptec controllers have the possibility to choose between the usual `H=64, S=32` and the 'extended translation' `H=255, S=63`.) Sometimes the BIOS reads the partition table to see with what geometry the disk was last partitioned - it will assume that a valid partition table is present when the `55aa` signature is present. This is good, in that it allows moving disks to a different machine. But having the BIOS behaviour depend on the disk contents also causes strange problems. (For example, it has been reported that a 2.5 GB disk was seen as having 528 MB because the BIOS read the partition table and concluded that it should use untranslated CHS. Another effect is found in the report that unpartitioned disks were slower than partitioned ones, because the BIOS tested 32-bit mode by reading the MBR and seeing whether it correctly got the `55aa` signature.)

How does the disk know about the geometry? Well, the manufacturer invents a geometry that multiplies out to approximately the right capacity. Many disks have jumpers that change the reported geometry, in order to avoid BIOS bugs. For example, all IBM disks allow the user to choose between 15 and 16 heads, and many manufacturers add jumpers to make the disk seem smaller than 2.1 GB or 33.8 GB. See also [below](#). Sometimes there are utilities that change the disk firmware.

9.1 Computing LILO parameters

Sometimes it is useful to force a certain geometry by adding ``hda=cyls,heads,secs'` on the kernel command line. Almost always one wants `secs=63`, and the purpose of adding this is to specify `heads`. (Reasonable values today are `heads=16` and `heads=255`.) What should one specify for `cyls`? Precisely that number that will give the right total capacity of $C*H*S$ sectors. For example, for a drive with 71346240 sectors (36529274880 bytes) one would compute C as $71346240/(255*63)=4441$ (for example using the program `bc`), and give boot parameter `hdc=4441,255,63`. How does one know the right total capacity? For example,

```
# hdparm -g /dev/hdc | grep sectors
geometry      = 4441/255/63, sectors = 71346240, start = 0
# hdparm -i /dev/hdc | grep LBAsects
CurCHS=16383/16/63, CurSects=16514064, LBA=yes, LBAsects=71346240
```

gives two ways of finding the total number of sectors 71346240. Recent kernels also give the precise size in the boot messages:

```
# dmesg | grep hde
hde: Maxtor 93652U8, ATA DISK drive
hde: 71346240 sectors (36529 MB) w/2048KiB Cache, CHS=70780/16/63
hde: hde1 hde2 hde3 < hde5 > hde4
hde2: <bsd: hde6 hde7 hde8 hde9 >
```

Older kernels only give MB and CHS. In general the CHS value is rounded down, so that the above output tells us that there are at least $70780*16*63=71346240$ sectors. In this example that happens to be the precise value. The MB value may be rounded instead of truncated, and in old kernels may be ``binary'` (MiB) instead of decimal. Note the agreement between the kernel size in MB and the Maxtor model number. Also in the case of SCSI disks the precise number of sectors is given in the kernel boot messages:

```
SCSI device sda: 17755792 512-byte hdwr sectors (9091 MB)
```

10. Details

10.1 IDE details - the seven geometries

The IDE driver has five sources of information about the geometry. The first (`G_user`) is the one specified by the user on the command line. The second (`G_bios`) is the BIOS Fixed Disk Parameter Table (for first and second disk only) that is read on system startup, before the switch to 32-bit mode. The third (`G_phys`) and fourth (`G_log`) are returned by the IDE controller as a response to the IDENTIFY command - they are the ``physical'` and ``current logical'` geometries.

On the other hand, the driver needs two values for the geometry: on the one hand `G_fdisk`, returned by a `HDIO_GETGEO` ioctl, and on the other hand `G_used`, which is actually used for doing I/O. Both `G_fdisk` and `G_used` are initialized to `G_user` if given, to `G_bios` when this information is present according to CMOS, and to `G_phys` otherwise. If `G_log` looks reasonable then `G_used` is set to that. Otherwise, if `G_used` is

Large Disk HOWTO

unreasonable and `G_phys` looks reasonable then `G_used` is set to `G_phys`. Here 'reasonable' means that the number of heads is in the range 1-16.

To say this in other words: the command line overrides the BIOS, and will determine what `fdisk` sees, but if it specifies a translated geometry (with more than 16 heads), then for kernel I/O it will be overridden by output of the IDENTIFY command.

Note that `G_bios` is rather unreliable: for systems booting from SCSI the first and second disk may well be SCSI disks, and the geometry that the BIOS reported for `sda` is used by the kernel for `hda`. Moreover, disks that are not mentioned in the BIOS Setup are not seen by the BIOS. This means that, e.g., in an IDE-only system where `hdb` is not given in the Setup, the geometries reported by the BIOS for the first and second disk will apply to `hda` and `hdc`.

In order to avoid such confusion, since Linux 2.5.51 `G_bios` is not used anymore.

The IDENTIFY DRIVE command

When an IDE drive is sent the IDENTIFY DRIVE (0xec) command, it will return 256 words (512 bytes) of information. This contains lots of technical stuff. Let us only describe here what plays a role in geometry matters. The words are numbered 0-255.

We find four pieces of information here: DefaultCHS (words 1,3,6), CurrentCHS (words 54-58) and LBACapacity (words 60-61) and 48-bit capacity (words 100-103).

	Description	Example
0	bit field: bit 6: fixed disk, bit 7: removable medium	0x0040
1	Default number of cylinders	16383
3	Default number of heads	16
6	Default number of sectors per track	63
10-19	Serial number (in ASCII)	G8067TME
23-26	Firmware revision (in ASCII)	GAK&1B0
27-46	Model name (in ASCII)	Maxtor 4G160J8
49	bit field: bit 9: LBA supported	0x2f00
53	bit field: bit 0: words 54-58 are valid	0x0007
54	Current number of cylinders	16383
55	Current number of heads	16
56	Current number of sectors per track	63
57-58	Current LBA capacity	16514064
60-61	Default LBA capacity	268435455
82-83	Command sets supported	7c69 4f09
85-86	Command sets enabled	7c68 0e01
100-103	Maximum user LBA for 48-bit addressing	320173056

Large Disk HOWTO

255	Checksum and signature (0xa5)	0x44a5
-----	-------------------------------	--------

In the ASCII strings each word contains two characters, the high order byte the first, the low order byte the second. The 32-bit values are given with low order word first. Words 54-58 are set by the command INITIALIZE DRIVE PARAMETERS (0x91). They are significant only when CHS addressing is used, but may help to find the actual disk size in case the disk sets DefaultCHS to 4092/16/63 in order to avoid BIOS problems.

Sometimes, when a jumper causes a big drive to misreport LBAcapacity (often to 66055248 sectors, in order to stay below the 33.8 GB limit), one needs a fourth piece of information to find the actual disk size, namely the result of the READ NATIVE MAX ADDRESS (0xf8) command.

10.2 SCSI details

The situation for SCSI is slightly different, as the SCSI commands already use logical block numbers, so a 'geometry' is entirely irrelevant for actual I/O. However, the format of the partition table is still the same, so `fdisk` has to invent some geometry, and also uses `HDIO_GETGEO` here - indeed, `fdisk` does not distinguish between IDE and SCSI disks. As one can see from the detailed description below, the various drivers each invent a somewhat different geometry. Indeed, one big mess.

If you are not using DOS or so, then avoid all extended translation settings, and just use 64 heads, 32 sectors per track (for a nice, convenient 1 MiB per cylinder), if possible, so that no problems arise when you move the disk from one controller to another. Some SCSI disk drivers (`aha152x`, `pas16`, `ppa`, `qlogicfas`, `qlogicisp`) are so nervous about DOS compatibility that they will not allow a Linux-only system to use more than about 8 GiB. This is a bug.

What is the real geometry? The easiest answer is that there is no such thing. And if there were, you wouldn't want to know, and certainly NEVER, EVER tell `fdisk` or `LILLO` or the kernel about it. It is strictly a business between the SCSI controller and the disk. Let me repeat that: only silly people tell `fdisk`/`LILLO`/kernel about the true SCSI disk geometry.

But if you are curious and insist, you might ask the disk itself. There is the important command `READ CAPACITY` that will give the total size of the disk, and there is the `MODE SENSE` command, that in the Rigid Disk Drive Geometry Page (page 04) gives the number of cylinders and heads (this is information that cannot be changed), and in the Format Page (page 03) gives the number of bytes per sector, and sectors per track. This latter number is typically dependent upon the notch, and the number of sectors per track varies - the outer tracks have more sectors than the inner tracks. The Linux program `scsiinfo` will give this information. There are many details and complications, and it is clear that nobody (probably not even the operating system) wants to use this information. Moreover, as long as we are only concerned about `fdisk` and `LILLO`, one typically gets answers like `C/H/S=4476/27/171` - values that cannot be used by `fdisk` because the partition table reserves only 10 resp. 8 resp. 6 bits for C/H/S.

Then where does the kernel `HDIO_GETGEO` get its information from? Well, either from the SCSI controller, or by making an educated guess. Some drivers seem to think that we want to know 'reality', but of course we only want to know what the DOS or OS/2 `FDISK` (or Adaptec `AFDISK`, etc) will use.

Large Disk HOWTO

Note that Linux `fdisk` needs the numbers H and S of heads and sectors per track to convert LBA sector numbers into c/h/s addresses, but the number C of cylinders does not play a role in this conversion. Some drivers use $(C,H,S) = (1023,255,63)$ to signal that the drive capacity is at least $1023 \times 255 \times 63$ sectors. This is unfortunate, since it does not reveal the actual size, and will limit the users of most `fdisk` versions to about 8 GiB of their disks - a real limitation in these days.

In the description below, M denotes the total disk capacity, and C, H, S the number of cylinders, heads and sectors per track. It suffices to give H, S if we regard C as defined by $M / (H \times S)$.

By default, H=64, S=32.

aha1740, dtc, g_NCR5380, t128, wd7000:

H=64, S=32.

aha152x, pas16, ppa, qllogicfas, qllogicisp:

H=64, S=32 unless $C > 1024$, in which case H=255, S=63, $C = \min(1023, M/(H \times S))$. (Thus C is truncated, and $H \times S \times C$ is not an approximation to the disk capacity M. This will confuse most versions of `fdisk`.) The `ppa.c` code uses M+1 instead of M and says that due to a bug in `sd.c` M is off by 1.

advansys:

H=64, S=32 unless $C > 1024$ and moreover the '> 1 GB' option in the BIOS is enabled, in which case H=255, S=63.

aha1542:

Ask the controller which of two possible translation schemes is in use, and use either H=255, S=63 or H=64, S=32. In the former case there is a boot message "aha1542.c: Using extended bios translation".

aic7xxx:

H=64, S=32 unless $C > 1024$, and moreover either the "extended" boot parameter was given, or the 'extended' bit was set in the SEEPROM or BIOS, in which case H=255, S=63. In Linux 2.0.36 this extended translation would always be set in case no SEEPROM was found, but in Linux 2.2.6 if no SEEPROM is found extended translation is set only when the user asked for it using this boot parameter (while when a SEEPROM is found, the boot parameter is ignored). This means that a setup that works under 2.0.36 may fail to boot with 2.2.6 (and require the `linear` keyword for LILO, or the `aic7xxx=extended` kernel boot parameter).

buslogic:

H=64, S=32 unless $C \geq 1024$, and moreover extended translation was enabled on the controller, in which case if $M < 2^{22}$ then H=128, S=32; otherwise H=255, S=63. However, after making this choice for (C,H,S), the partition table is read, and if for one of the three possibilities (H,S) = (64,32), (128,32), (255,63) the value `endH=H-1` is seen somewhere then that pair (H,S) is used, and a boot message is printed "Adopting Geometry from Partition Table".

fdomain:

Find the geometry information in the BIOS Drive Parameter Table, or read the partition table and use `H=endH+1, S=endS` for the first partition, provided it is nonempty, or use H=64, S=32 for $M < 2^{21}$ (1 GiB), H=128, S=63 for $M < 63 \times 2^{17}$ (3.9 GiB) and H=255, S=63 otherwise.

in2000:

Use the first of (H,S) = (64,32), (64,63), (128,63), (255,63) that will make $C \leq 1024$. In the last case, truncate C at 1023.

seagate:

Read C,H,S from the disk. (Horrors!) If C or S is too large, then put S=17, H=2 and double H until $C \leq 1024$. This means that H will be set to 0 if $M > 128 \times 1024 \times 17$ (1.1 GiB). This is a bug.

ultrastor and u14_34f:

One of three mappings ((H,S) = (16,63), (64,32), (64,63)) is used depending on the controller mapping mode.

Large Disk HOWTO

If the driver does not specify the geometry, we fall back on an educated guess using the partition table, or using the total disk capacity.

Look at the partition table. Since by convention partitions end on a cylinder boundary, we can, given $end = (endC, endH, endS)$ for any partition, just put $H = endH + 1$ and $S = endS$. (Recall that sectors are counted from 1.) More precisely, the following is done. If there is a nonempty partition, pick the partition with the largest $beginC$. For that partition, look at $end + 1$, computed both by adding $start$ and $length$ and by assuming that this partition ends on a cylinder boundary. If both values agree, or if $endC = 1023$ and $start + length$ is an integral multiple of $(endH + 1) * endS$, then assume that this partition really was aligned on a cylinder boundary, and put $H = endH + 1$ and $S = endS$. If this fails, either because there are no partitions, or because they have strange sizes, then look only at the disk capacity M . Algorithm: put $H = M / (62 * 1024)$ (rounded up), $S = M / (1024 * H)$ (rounded up), $C = M / (H * S)$ (rounded down). This has the effect of producing a (C, H, S) with C at most 1024 and S at most 62.

11. Clipped disks

11.1 The Linux IDE 8 GiB limit

The Linux IDE driver gets the geometry and capacity of a disk (and lots of other stuff) by using an [ATA IDENTIFY](#) request. Linux kernels older than 2.0.34/2.1.90 would not believe the returned value of `lba_capacity` if it was more than 10% larger than the capacity computed by $C * H * S$. However, by industry agreement large IDE disks (with more than 16514064 sectors) return $C=16383$, $H=16$, $S=63$, for a total of 16514064 sectors (7.8 GB) independent of their actual size, but give their actual size in `lba_capacity`.

Since versions 2.0.34/2.1.90, Linux kernels know about this and do the right thing. If you have an older Linux kernel and do not want to upgrade, and this kernel only sees 8 GiB of a much larger disk, then try changing the routine `lba_capacity_is_ok` in `/usr/src/linux/drivers/block/ide.c` into something like

```
static int lba_capacity_is_ok (struct hd_driveid *id) {
    id->cylns = id->lba_capacity / (id->heads * id->sectors);
    return 1;
}
```

For a more cautious patch, see 2.1.90.

11.2 BIOS complications

As just mentioned, large disks return the geometry $C=16383$, $H=16$, $S=63$ independent of the actual size, while the actual size is returned in the value of `LBACapacity`. Some BIOSes do not recognize this, and translate this 16383/16/63 into something with fewer cylinders and more heads, for example 1024/255/63 or 1027/255/63. So, the kernel must not only recognize the single geometry 16383/16/63, but also all BIOS-mangled versions of it. Since 2.2.2 this is done correctly (by taking the BIOS idea of H and S , and computing $C = capacity / (H * S)$). Usually this problem is solved by setting the disk to Normal in the BIOS setup (or, even better, to None, not mentioning it at all to the BIOS). If that is impossible because you have to boot from it or use it also with DOS/Windows, and upgrading to 2.2.2 or later is not an option, use kernel boot parameters.

If a BIOS reports 16320/16/63, then this is usually done in order to get 1024/255/63 after translation.

Large Disk HOWTO

There is an additional problem here. If the disk was partitioned using a geometry translation, then the kernel may at boot time see this geometry used in the partition table, and report `hda: [PTBL] [1027/255/63]`. This is bad, because now the disk is only 8.4 GB. This was fixed in 2.3.21. Again, kernel boot parameters will help.

11.3 Jumpers that select the number of heads

Many disks have jumpers that allow you to choose between a 15-head and a 16-head geometry. The default settings will give you a 16-head disk. Sometimes both geometries address the same number of sectors, sometimes the 15-head version is smaller. There may be a good reason for this setup: Petri Kaukasoina writes: 'A 10.1 Gig IBM Deskstar 16 GP (model IBM-DTTA-351010) was jumpered for 16 heads as default but this old PC (with AMI BIOS) didn't boot and I had to jumper it for 15 heads. `hdparm -i` tells `RawCHS=16383/15/63` and `LBAssects=19807200`. I use `20960/15/63` to get the full capacity.' For the jumper settings, see <http://www.hitachigst.com/hdd/support/jumpers.htm>.

11.4 Jumpers that clip total capacity

Many disks have jumpers that allow you to make the disk appear smaller than it is. A silly thing to do, and probably no Linux user ever wants to use this, but some BIOSes crash on big disks. The usual solution is to keep the disk entirely out of the BIOS setup. But this may be feasible only if the disk is not your boot disk.

Clip to 2.1 GB

The first serious limit was the 4096 cylinder limit (that is, with 16 heads and 63 sectors/track, 2.11 GB). For example, a Fujitsu MPB3032ATU 3.24 GB disk has default geometry 6704/15/63, but can be jumpered to appear as 4092/16/63, and then reports `LBACapacity 4124736` sectors, so that the operating system cannot guess that it is larger in reality. In such a case (with a BIOS that crashes if it hears how big the disk is in reality, so that the jumper is required) one needs boot parameters to tell Linux about the size of the disk.

That is unfortunate. Most disks can be jumpered so as to appear as a 2 GB disk and then report a clipped geometry like 4092/16/63 or 4096/16/63, but still report full `LBACapacity`. Such disks will work well, and use full capacity under Linux, regardless of jumper settings.

Clip to 33 GB

A more recent limit is [the 33.8 GB limit](#). Linux kernels older than 2.2.14 / 2.3.21 need a patch to be able to cope with IDE disks larger than this.

With an old BIOS and a disk larger than 33.8 GB, the BIOS may hang, and in such cases booting may be impossible, even when the disk is removed from the CMOS settings.

Therefore, large IBM and Maxtor and Seagate disks come with a jumper that make the disk appear as a 33.8 GB disk. For example, the IBM Deskstar 37.5 GB (DPTA-353750) with 73261440 sectors (corresponding to 72680/16/63, or 4560/255/63) can be jumpered to appear as a 33.8 GB disk, and then reports geometry 16383/16/63 like any big disk, but `LBACapacity 66055248` (corresponding to 65531/16/63, or 4111/255/63). Similar things hold for recent large Maxtor disks.

Below some more details that used to be relevant but probably can be ignored now.

Maxtor

With the jumper present, both the geometry (16383/16/63) and the size (66055248) are conventional and give no information about the actual size. Moreover, attempts to access sector 66055248 and above yield I/O errors. However, on Maxtor drives the actual size can be found and made accessible using the READ NATIVE MAX ADDRESS and SET MAX ADDRESS commands. Presumably this is what MaxBlast/EZ-Drive does. There is a small Linux utility [setmax.c](#) that does the same. Only very few disks need it - almost always CONFIG_IDEDISK_STROKE does the trick.

For drives larger than 137 GB also READ NATIVE MAX ADDRESS returns a conventional value, namely 0xffffffff, corresponding to 137 GB. Here READ NATIVE MAX ADDRESS EXT and SET MAX ADDRESS EXT (using 48-bit addressing) are required. The `setmax` utility does not yet know about this. A very small patch makes 2.5.3 handle this situation.

Early large Maxtor disks have an additional detail: the J46 jumper for these 34-40 GB disks changes the geometry from 16383/16/63 to 4092/16/63 and does not change the reported LBA capacity. This means that also with jumper present the BIOS (old Award 4.5*) will hang at boot time. For this case Maxtor provides a utility [JUMPON.EXE](#) that upgrades the firmware to make J46 behave as described above.

On recent Maxtor drives the call `setmax -d 0 /dev/hdX` will give you max capacity again. However, on slightly older drives a firmware bug does not allow you to use `-d 0`, and `setmax -d 255 /dev/hdX` returns you to almost full capacity. For Maxtor D540X-4K, see below.

IBM

For IBM things are worse: the jumper really clips capacity and there is no software way to get it back. The solution is not to use the jumper but use `setmax -m 66055248 /dev/hdX` to software-clip the disk. "How?" you say - "I cannot boot!". IBM gives the tip: *If a system with Award BIOS hangs during drive detection: Reboot the system and hold the F4 key to bypass autodetection of the drive(s).* If this doesn't help, find a different computer, connect the drive to it, and run `setmax` there. After doing this you go back to the first machine and are in the same situation as with jumpered Maxtor disks: booting works, and after getting past the BIOS either a patched kernel or a `setmax -d 0` gets you full capacity.

Thomas Charbonnel reports on a different approach: "I had a 80 GB IBM IC35L080AVVA07-0 drive and installed IBM's Disk Manager. Installed my boot loader on the drive's MBR. Everything worked fine. Note that the IDE drive must become the boot drive so that one can install only one 34+ GB drive using this approach."

Seagate

Seagate disks have a jumper that will clip the reported number of cylinders to 4092 on drives smaller than 33.8 GB, while it will limit the reported LBA capacity (Identify words 60/61) to 33.8 GB on larger disks.

For models ST-340810A, ST-360020A, ST-380020A: The ATA Read Native Max and Set Max commands may be used to reset the true full capacity.

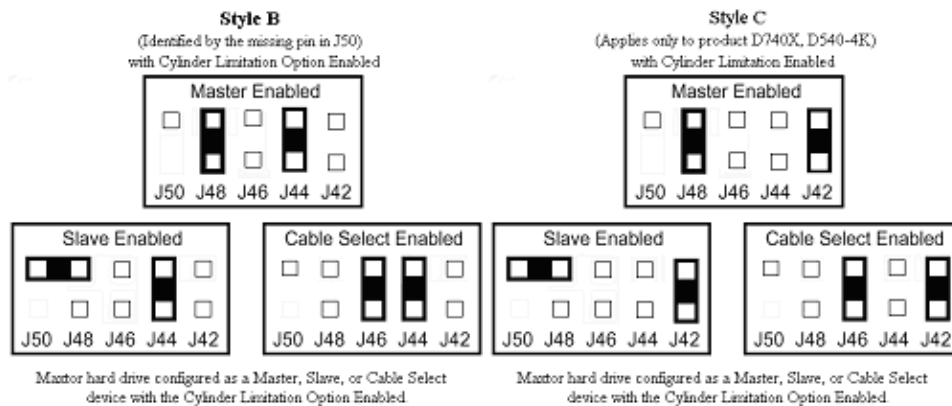
For models ST-340016A, ST-340823A, ST-340824A, ST-360021A, ST-380021A: The ATA Set Features F1 sub-command will cause Identify Data words 60-61 to report the true full capacity.

Maxtor D540X-4K

The Maxtor Diamond Max drives 4K080H4, 4K060H3, 4K040H2 (aka D540X-4K) are identical to the drives 4D080H4, 4D060H3, 4D040H2 (aka D540X-4D), except that the jumper settings differ. A Maxtor FAQ specifies the Master/Slave/CableSelect settings for them, but the capacity clip jumper for the "4K" drives seems to be undocumented. Nils Ohlmeier reports that he experimentally finds that it is the J42 jumper ("reserved for factory use") closest to the power connector. (The "4D" drives use the J46 jumper, like all other Maxtor drives.)

However, it may be that this undocumented jumper acts like the IBM jumper: the machine boots correctly, but the disk has been clipped to 33 GB and `setmax -d 0` does not help to get full capacity back. And the IBM solution works: do not use any disk-clipping jumpers, but first put the disk in a machine with non-broken BIOS, soft-clip it with `setmax -m 66055248 /dev/hdX`, then put it back in the first machine, and after booting run `setmax -d 0 /dev/hdX` to get full capacity again.

In the meantime, some docs and pictures have appeared on the Maxtor site, confirming part of the above. Compare



Western Digital

Some info, including the settings for capacity-clipping jumpers, is given on [the Western Digital site](#). I do not know what precisely these jumpers do.

11.5 READ NATIVE MAX ADDRESS / SET MAX ADDRESS

If an IDE/ATA disk has support for the Host Protected Area (HPA) feature set, then it is possible to set the LBA capacity to any value below the actual capacity. Access past the assigned point usually leads to I/O errors. Since classical software finds out about the disk size by looking at the LBA capacity field of the Identify information, such software will not suspect that the disk actually is larger.

The actual total size of the disk is read using the READ NATIVE MAX ADDRESS command. This "soft disk size" is set using the SET MAX ADDRESS command. It comes in two flavours: if the "volatile" bit is set, the command will have effect until the next reboot or hardware reset; otherwise the effect is permanent. It is possible to protect settings with a password. (For details, see the ATA standard.)

This clipped size has (at least) two applications: on the one hand it is possible to fake a smaller disk, so that the BIOS will not have problems, and have Linux, or (for DOS/Windows) a disk manager restore total size;

on the other hand one can have a vendor area at the end, inaccessible to the ordinary user.

For many of the disks discussed above, setting a jumper has precisely this effect: LBA capacity is diminished while the native max capacity remains the same, and the SET MAX ADDRESS will restore full capacity.

11.6 CONFIG_IDEDISK_STROKE

The CONFIG_IDEDISK_STROKE option of Linux 2.4.19/2.5.3 and later, will tell Linux to read the native max capacity and do a SET MAX ADDRESS to get access to full capacity. This configuration option lives under the heading "Auto-Geometry Resizing support" in the "IDE, ATA and ATAPI block devices" kernel configuration section.

The configuration option went away in 2.6.7 and was replaced by a (per-disk) boot parameter, so that one can say "hda=stroke".

With this "stroke" option jumpered disks will in many cases be handled correctly, i.e., be seen with full capacity (in spite of the jumper). And the same holds when the disk got a Host Protected Area in some other (non-jumper) way.

This is the preferred way to handle disks that need a jumper because of a broken BIOS.

12. The Linux 65535 cylinder limit

The HDIO_GETGEO ioctl returns the number of cylinders in a short. This means that if you have more than 65535 cylinders, the number is truncated, and (for a typical SCSI setup with 1 MiB cylinders) a 80 GiB disk may appear as a 16 GiB one. Once one recognizes what the problem is, it is easily avoided. Use fdisk 2.10i or newer.

(The programming convention is to use the BLKGETSIZE ioctl to get total size, and HDIO_GETGEO to get number of heads and sectors/track, and, if needed, get C by $C = \text{size}/(H*S)$.)

12.1 IDE problems with 34+ GB disks

(Below a discussion of Linux kernel problems. BIOS problems and jumpers that clip capacity were discussed [above](#).)

Drives larger than 33.8 GB will not work with kernels older than 2.0.39 / 2.2.14 / 2.3.21. The details are as follows. Suppose you bought a new IBM-DPTA-373420 disk with a capacity of 66835440 sectors (34.2 GB). Pre-2.3.21 kernels will tell you that the size is $769*16*63 = 775152$ sectors (0.4 GB), which is a bit disappointing. And giving command line parameters `hdc=4160,255,63` doesn't help at all - these are just ignored. What happens? The routine `idedisk_setup()` retrieves the geometry reported by the disk (which is 16383/16/63) and overwrites what the user specified on the command line, so that the user data is used only for the BIOS geometry. The routine `current_capacity()` or `idedisk_capacity()` recomputes the cylinder number as $66835440/(16*63)=66305$, but since this is stored in a short, it becomes 769. Since `lba_capacity_is_ok()` destroyed `id->cyls`, every following call to it will return false, so that the disk capacity becomes $769*16*63$. For several kernels a patch is available. A patch for 2.0.38 can be found at ftp.kernel.org. A patch for 2.2.12 can be found at www.uwsg.indiana.edu (some editing may be required to get rid of the html markup). The 2.2.14 kernels do support these disks. In the 2.3.* kernel series, there is support for these disks since 2.3.21. One can also 'solve' the problem in hardware by [using a jumper](#) to clip the size to 33.8 GB. In many cases a

BIOS upgrade will be required if one wants to boot from the disk.

13. Extended and logical partitions

Above, we saw the structure of the MBR (sector 0): boot loader code followed by 4 partition table entries of 16 bytes each, followed by an AA55 signature. Partition table entries of type 5 or F or 85 (hex) have a special significance: they describe *extended* partitions: blobs of space that are further partitioned into *logical* partitions. (So, an extended partition is only a box, it cannot be used itself, one uses the logical partitions inside.) Only the location of the first sector of an extended partition is important. This first sector contains a partition table with four entries: one a logical partition, one an extended partition, and two unused. In this way one gets a chain of partition table sectors, scattered over the disk, where the first one describes three primary partitions and the extended partition, and each following partition table sector describes one logical partition and the location of the next partition table sector.

It is important to understand this: When people do something stupid while partitioning a disk, they want to know: Is my data still there? And the answer is usually: Yes. But if logical partitions were created then the partition table sectors describing them are written at the beginning of these logical partitions, and data that was there before is lost.

The program `sfdisk` will show the full chain. E.g.,

```
# sfdisk -l -x /dev/hda

Disk /dev/hda: 16 heads, 63 sectors, 33483 cylinders
Units = cylinders of 516096 bytes, blocks of 1024 bytes, counting from 0

   Device Boot  Start      End  #cyls   #blocks  Id System
/dev/hda1            0+    101    102-    51376+  83 Linux
/dev/hda2           102   2133   2032   1024128  83 Linux
/dev/hda3           2134  33482  31349  15799896   5 Extended
/dev/hda4            0        -     0         0    0 Empty

/dev/hda5           2134+   6197   4064-   2048224+  83 Linux
-                   6198  10261   4064   2048256   5 Extended
-                   2134   2133     0         0    0 Empty
-                   2134   2133     0         0    0 Empty

/dev/hda6           6198+  10261   4064-   2048224+  83 Linux
-                   10262 16357   6096   3072384   5 Extended
-                   6198   6197     0         0    0 Empty
-                   6198   6197     0         0    0 Empty
...
/dev/hda10          30581+ 33482  2902-  1462576+  83 Linux
-                   30581 30580     0         0    0 Empty
-                   30581 30580     0         0    0 Empty
-                   30581 30580     0         0    0 Empty

#
```

It is possible to construct bad partition tables. Many kernels get into a loop if some extended partition points back to itself or to an earlier partition in the chain. It is possible to have two extended partitions in one of these partition table sectors so that the partition table chain forks. (This can happen for example with an `fdisk` that does not recognize each of 5, F, 85 as an extended partition, and creates a 5 next to an F.) No standard `fdisk` type program can handle such situations, and some handwork is required to repair them. The Linux kernel will accept a fork at the outermost level. That is, you can have two chains of logical partitions.

Large Disk HOWTO

Sometimes this is useful - for example, one can use type 5 and be seen by DOS, and the other type 85, invisible for DOS, so that DOS FDISK will not crash because of logical partitions past cylinder 1024. Usually one needs `sfdisk` to create such a setup.

14. Problem solving

Many people think they have problems, while in fact nothing is wrong. Or, they think that the problems they have are due to disk geometry, while in fact disk geometry has nothing to do with the matter. All of the above may have sounded complicated, but disk geometry handling is extremely easy: do nothing at all, and all is fine; or perhaps give LILO the keyword `lba32` if it doesn't get past ``LI'` when booting. Watch the kernel boot messages, and remember: the more you fiddle with geometries (specifying heads and cylinders to LILO and `fdisk` and on the kernel command line) the less likely it is that things will work. Roughly speaking, all is fine by default.

And remember: nowhere in Linux is disk geometry used, so no problem you have while running Linux can be caused by disk geometry. Indeed, disk geometry is used only by LILO and by `fdisk`. So, if LILO fails to boot the kernel, that may be a geometry problem. If different operating systems do not understand the partition table, that may be a geometry problem. Nothing else. In particular, if `mount` doesn't seem to work, never worry about disk geometry - the problem is elsewhere.

14.1 Problem: My IDE disk gets a bad geometry when I boot from SCSI.

It is quite possible that a disk gets the wrong geometry. The Linux kernel asks the BIOS about `hd0` and `hd1` (the BIOS drives numbered 80H and 81H) and assumes that this data is for `hda` and `hdb`. But on a system that boots from SCSI, the first two disks may well be SCSI disks, and thus it may happen that the fifth disk, which is the first IDE disk `hda`, gets assigned a geometry belonging to `sda`. Such things are easily solved by giving boot parameters ``hda=C,H,S'` for the appropriate numbers C, H and S, either at boot time or in `/etc/lilo.conf`.

Since Linux 2.5.51 this BIOS information is not used anymore, and the same problem occurs for all disks. See below.

14.2 Nonproblem: Identical disks have different geometry?

``I` have two identical 10 GB IBM disks. However, `fdisk` gives different sizes for them. Look:

```
# fdisk -l /dev/hdb
Disk /dev/hdb: 255 heads, 63 sectors, 1232 cylinders
Units = cylinders of 16065 * 512 bytes

   Device Boot   Start       End   Blocks   Id  System
/dev/hdb1            1       1232  9896008+  83  Linux native
# fdisk -l /dev/hdd
Disk /dev/hdd: 16 heads, 63 sectors, 19650 cylinders
Units = cylinders of 1008 * 512 bytes

   Device Boot   Start       End   Blocks   Id  System
/dev/hdd1            1      19650  9903568+  83  Linux native
```

How come?'

Large Disk HOWTO

What is happening here? Well, first of all these drives really are 10gig: hdb has size $255 * 63 * 1232 * 512 = 10133544960$, and hdd has size $16 * 63 * 19650 * 512 = 10141286400$, so, nothing is wrong and the kernel sees both as 10.1 GB. Why the difference in size? That is because the kernel gets data for the first two IDE disks from the BIOS, and the BIOS has remapped hdb to have 255 heads (and $16 * 19650 / 255 = 1232$ cylinders). The rounding down here costs almost 8 MB.

If you would like to remap hdd in the same way, give the kernel boot parameters ``hdd=1232,255,63'`.

On the other hand, if the disk is not shared with DOS or so, it may be better to set hdb to Normal in the BIOS setup, instead of asking for some translation like LBA.

Since Linux 2.5.51, the IDE driver no longer uses BIOS info on the first two disks, and the different treatment of the first two disks has disappeared.

14.3 Problem: 2.4 and 2.6 report different geometries? 2.6 reports the wrong geometry? 2.6 reports no geometry at all?

Since geometry does not exist, it is not surprising that each of 2.0/2.2/2.4/2.6 reports a somewhat different disk geometry.

Some people will maintain that geometry **does** exist, and in that case do not mean a property of the disk, but mean the values reported by the BIOS. That is what several other operating systems will use. Since Linux 2.5.51, the kernel no longer uses the values reported by the BIOS - it is difficult to match BIOS device numbers with Linux disk names, maybe data is only available for two disks, maybe some disks are not present in the BIOS setup, etc. However, if one needs these values, since Linux 2.6.5 one can set CONFIG_EDD and mount sysfs, and then find the BIOS data for the various disks under `/sys/firmware/edd/int13_dev*`. Now the matching of BIOS numbers, represented in directory names like `int13_dev82`, with Linux names like `sda` can be done by user space software, possibly with help from the user.

This 2.5.51 change caused problems when many people using both Linux and Windows on the same disk upgraded from 2.4 to 2.6 and used as partitioning tool the program `parted` that had not yet been updated. I have not checked whether current `parted` is OK.

14.4 Nonproblem: fdisk sees much more room than df?

`fdisk` will tell you how many blocks there are on the disk. If you make a filesystem on the disk, say with `mke2fs`, then this filesystem needs some space for bookkeeping - typically something like 4% of the filesystem size, more if you ask for a lot of inodes during `mke2fs`. For example:

```
# sfdisk -s /dev/hda9
4095976
# mke2fs -i 1024 /dev/hda9
mke2fs 1.12, 9-Jul-98 for EXT2 FS 0.5b, 95/08/09
...
204798 blocks (5.00%) reserved for the super user
...
# mount /dev/hda9 /somewhere
# df /somewhere
```

Large Disk HOWTO

```
Filesystem      1024-blocks  Used Available Capacity Mounted on
/dev/hda9       3574475      13 3369664      0% /mnt
# df -i /somewhere
Filesystem      Inodes      IUsed   IFree  %IUsed Mounted on
/dev/hda9       4096000     11 4095989      0% /mnt
#
```

We have a partition with 4095976 blocks, make an ext2 filesystem on it, mount it somewhere and find that it only has 3574475 blocks - 521501 blocks (12%) was lost to inodes and other bookkeeping. Note that the difference between the total 3574475 and the 3369664 available to the user are the 13 blocks in use plus the 204798 blocks reserved for root. This latter number can be changed by `tune2fs`. This `-i 1024` is only reasonable for news spools and the like, with lots and lots of small files. The default would be:

```
# mke2fs /dev/hda9
# mount /dev/hda9 /somewhere
# df /somewhere
Filesystem      1024-blocks  Used Available Capacity Mounted on
/dev/hda9       3958475      13 3753664      0% /mnt
# df -i /somewhere
Filesystem      Inodes      IUsed   IFree  %IUsed Mounted on
/dev/hda9       1024000     11 1023989      0% /mnt
#
```

Now only 137501 blocks (3.3%) are used for inodes, so that we have 384 MB more than before. (Apparently, each inode takes 128 bytes.) On the other hand, this filesystem can have at most 1024000 files (more than enough), against 4096000 (too much) earlier.