

What Partitions do I need?	7
Discussion:	8
File Systems	8
Swap Partitions	9
Partitioning with fdisk	11
fdisk usage	11
Four primary partitions	12
Mixed primary and logical partitions	13
Submitted Examples	15
Labels	15
Volume Labels	16
Device Labels	16
Formatting an ext2/3 partition	17
Simple Invocation	17
Reserved blocks	17
Recovering a Deleted Partition Table	17
Setting Up Swap Space	19
Swap Files	19
Swap Files	19
Multiple Swap Areas	20
Appendix	20
Formating Partitions	20
Activating Swap Space	21
Mounting Partitions	21
Some facts about file systems and fragmentation	21

Introduction

What is a partition?

Partitioning is a means to divide a single hard drive into many logical drives. A partition is a contiguous set of blocks on a drive that are treated as an independant disk. A partition table (the creation of which is the topic of this HOWTO) is an index that relates sections of the hard drive to partitions.

Why have multiple partitions?

- Encapsulate your data. Since file system corruption is local to a partition, you stand to lose only some of your data if an accident occurs.
- Increase disk space efficiency. You can format partitions with varying block sizes, depending on your usage. If your data is in a large number of small files (less than 1k) and your partition uses 4k sized blocks, you are wasting 3k for every file. In general, you waste on average one half of a block for every file, so matching block size to the average size of your files is important if you have many files.
- Limit data growth. Runaway processes or maniacal users can consume so much disk space that the operating system no longer has room on the hard drive for its bookkeeping operations. This will lead to disaster. By segregating space, you ensure that things other than the operating system die when allocated disk space is exhausted.

Other Partitioning Software:

- **sfdisk**: a command-line version of fdisk

- **cfdisk**: a curses-based version of fdisk
- **parted** [<http://www.gnu.org/software/parted/parted.html>]: Gnu partition editor
- Partition Magic™ [<http://www.powerquest.com/partitionmagic/index.html>]: a commercial utility to create, resize, merge and convert partitions, without destroying data.
- Disk Drake™ [<http://www.linux-mandrake.com/diskdrake>]: a Perl/Gtk program to create, resize, and delete partitions

Related HOWTOs

Table 1. Related HOWTOs

Title	Author	Description
Dual boot install strategies [http://www.tldp.org/HOWTO/Install-Strategies]	Gjoen Stein [mailto:tobbyb@usa.net]	How to estimate the various size and speed requirements for different parts of the filesystem.
Linux Multiple Disk System Tuning [http://www.nyx.net/~sgjoen/disk.html]	Gjoen Stein [mailto:sgjoen@mail.nyx.net]	How to estimate the various size and speed requirements for different parts of the filesystem.
Linux Large Disk [http://www.tldp.org/HOWTO/Large-Disk-HOWTO.html]	Andries Brouwer [mailto:aeb@cw.nl]	Instructions and considerations regarding disks with more than 1024 cylinders
Linux Quota [http://metalab.unc.edu/mdw/HOWTO/mini/Quota.html]	Ralf van Dooren [mailto:mailto:r.vdooren@snow.nl]	Instructions on limiting disk space usage per user (quotas)
Partition-Rescue mini-HOWTO [http://www.linuxdocs.org/HOWTOs/mini/Partition-Rescue.html]	Jean-Daniel Dodin [mailto:jdanield@dodin.net]	How to restore linux partitions after they have been deleted by a Windows install. Does not appear to preserve data.
Linux ADSM Backup [http://www.tldp.org/HOWTO/ADSM-Backup.html]	Thomas Koenig [mailto:Thomas.Koenig@ciw.uni-karlsruhe.de]	Instructions on integrating Linux into an IBM ADSM backup environment.
Linux Backup with MS-DOS [http://www.ibiblio.org/pub/Linux/docs/HOWTO/other-formats/html_single/Backup-With-MSDOS.html]	Christopher Neufeld [mailto:neufeld@physics.utoronto.ca]	Information about MS-DOS drivers for Linux backups.
Linux HOWTO Index	Tim Bynum [mailto:linux-howto@sunsite.unc.edu]	Instructions on writing and submitting a HOWTO document

Additional information on your system:

- `/usr/src/linux/Documentation` [file:/usr/src/linux/Documentation]
- `ide.txt` [file:/usr/src/linux/Documentation/ide.txt]: Info about your IDE drivers
- `scsi.txt` [file:/usr/src/linux/Documentation/scsi.txt]: Info about your SCSI drivers

Devices

There is a special nomenclature that linux uses to refer to hard drive partitions that must be understood in order to follow the discussion on the following pages.

In Linux, partitions are represented by device files. These are phoney files located in `/dev`. Here are a few entries:

```
brw-rw----    1 root    disk      3,   0 May  5 1998 hda
brw-rw----    1 root    disk      8,   0 May  5 1998 sda
crw-----    1 root    tty       4,  64 May  5 1998 ttyS0
```

A device file is a file with type `c` (for "character" devices, devices that do not use the buffer cache) or `b` (for "block" devices, which go through the buffer cache). In Linux, all disks are represented as block devices only.

Device names

Naming Convention

By convention, IDE drives will be given device names `/dev/hda` to `/dev/hdd`. *Hard Drive A* (`/dev/hda`) is the first drive and *Hard Drive C* (`/dev/hdc`) is the third.

Table 2. IDE controller naming convention

drive name	drive controller	drive number
<code>/dev/hda</code>	1	1
<code>/dev/hdb</code>	1	2
<code>/dev/hdc</code>	2	1
<code>/dev/hdd</code>	2	2

A typical PC has two IDE controllers, each of which can have two drives connected to it. For example, `/dev/hda` is the first drive (master) on the first IDE controller and `/dev/hdd` is the second (slave) drive on the second controller (the fourth IDE drive in the computer).

You can write to these devices directly (using `cat` or `dd`). However, since these devices represent the entire disk, starting at the first block, you can mistakenly overwrite the master boot record and the partition table, which will render the drive unusable.

Table 3. partition names

drive name	drive controller	drive number	partition type	partition number
<code>/dev/hda1</code>	1	1	primary	1
<code>/dev/hda2</code>	1	1	primary	2
<code>/dev/hda3</code>	1	1	primary	3
<code>/dev/hda4</code>	1	1	swap	NA
<code>/dev/hdb1</code>	1	2	primary	1
<code>/dev/hdb2</code>	1	2	primary	2
<code>/dev/hdb3</code>	1	2	primary	3

/dev/hdb4	1	2	primary	4
-----------	---	---	---------	---

Once a drive has been partitioned, the partitions will be represented as numbers on the end of the names. For example, the second partition on the second drive will be `/dev/hdb2`. The partition type (primary) is listed in the table above for clarity, although the concept is not explained until the section called “Primary Partitions”.

Table 4. SCSI Drives

drive name	drive controller	drive number	partition type	partition number
/dev/sda1	1	6	primary	1
/dev/sda2	1	6	primary	2
/dev/sda3	1	6	primary	3

SCSI drives follow a similar pattern; They are represented by 'sd' instead of 'hd'. The first partition of the second SCSI drive would therefore be `/dev/sdb1`. In the table above, the drive number is arbitrarily chosen to be 6 to introduce the idea that SCSI ID numbers do not map onto device names under linux.

Name Assignment

Under (Sun) Solaris and (SGI) IRIX, the device name given to a SCSI drive has some relationship to where you plug it in. Under linux, there is only wailing and gnashing of teeth.

Before

SCSI ID #2	SCSI ID #5	SCSI ID #7	SCSI ID #8
/dev/sda	/dev/sdb	/dev/sdc	/dev/sdd

After

SCSI ID #2	SCSI ID #7	SCSI ID #8
/dev/sda	/dev/sdb	/dev/sdc

SCSI drives have ID numbers which go from 1 through 15. Lower SCSI ID numbers are assigned lower-order letters. For example, if you have two drives numbered 2 and 5, then #2 will be `/dev/sda` and #5 will be `/dev/sdb`. If you remove either, all the higher numbered drives will be renamed the next time you boot up.

If you have two SCSI controllers in your linux box, you will need to examine the output of `/bin/dmesg` in order to see what name each drive was assigned. If you remove one of two controllers, the remaining controller might have all its drives renamed. Grrr...

There are two work-arounds; both involve using a program to put a label on each partition (see the section called “Labels”). The label is persistent even when the device is physically moved. You then refer to the partition directly or indirectly by label.

Logical Partitions

Table 5. Logical Partitions

drive name	drive controller	drive number	partition type	partition number
/dev/hdb1	1	2	primary	1

/dev/hdb2	1	2	extended	NA
/dev/hda5	1	2	logical	2
/dev/hdb6	1	2	logical	3

The table above illustrates a mysterious jump in the name assignments. This is due to the use of logical partitions (see the section called “Logical Partitions”, which always start with 5, for reasons explained later.

This is all you have to know to deal with linux disk devices. For the sake of completeness, see Kristian's discussion of device numbers below.

Device numbers

The only important thing with a device file are its major and minor device numbers, which are shown instead of the file size:

```
$ ls -l /dev/hda
```

Table 6. Device file attributes

brw-rw----	1	root	disk	3,	0	Jul 18 1994	/dev/hda
permissions		owner	group	major device number	minor device number	date	device name

When accessing a device file, the major number selects which device driver is being called to perform the input/output operation. This call is being done with the minor number as a parameter and it is entirely up to the driver how the minor number is being interpreted. The driver documentation usually describes how the driver uses minor numbers. For IDE disks, this documentation is in `/usr/src/linux/Documentation/ide.txt` [file:/usr/src/linux/Documentation/ide.txt]. For SCSI disks, one would expect such documentation in `/usr/src/linux/Documentation/scsi.txt` [file:/usr/src/linux/Documentation/scsi.txt], but it isn't there. One has to look at the driver source to be sure (`/usr/src/linux/driver/scsi/sd.c` [file:/usr/src/linux/driver/scsi/sd.c]:184-196). Fortunately, there is Peter Anvin's list of device numbers and names in `/usr/src/linux/Documentation/devices.txt` [file:/usr/src/linux/Documentation/devices.txt]; see the entries for block devices, major 3, 22, 33, 34 for IDE and major 8 for SCSI disks. The major and minor numbers are a byte each and that is why the number of partitions per disk is limited.

Partition Types

Partition Types

A partition is labeled to host a certain kind of file system (not to be confused with a volume label (see the section called “Labels”)). Such a file system could be the linux standard ext2 file system or linux swap space, or even foreign file systems like (Microsoft) NTFS or (Sun) UFS. There is a numerical code associated with each partition type. For example, the code for ext2 is 0x83 and linux swap is 0x82. To see a list of partition types and their codes, execute `/sbin/sfdisk -T`

Foreign Partition Types

The partition type codes have been arbitrarily chosen (you can't figure out what they should be) and they are particular to a given operating system. Therefore, it is theoretically possible that if you use two operating

systems with the same hard drive, the same code might be used to designate two different partition types. OS/2 marks its partitions with a 0x07 type and so does Windows NT's NTFS. MS-DOS allocates several type codes for its various flavors of FAT file systems: 0x01, 0x04 and 0x06 are known. DR-DOS used 0x81 to indicate protected FAT partitions, creating a type clash with Linux/Minix at that time, but neither Linux/Minix nor DR-DOS are widely used any more.

OS/2 marks its partitions with a 0x07 type and so does Windows NT's NTFS. MS-DOS allocates several type codes for its various flavors of FAT file systems: 0x01, 0x04 and 0x06 are known. DR-DOS used 0x81 to indicate protected FAT partitions, creating a type clash with Linux/Minix at that time, but neither Linux/Minix nor DR-DOS are widely used any more.

Primary Partitions

The number of partitions on an Intel-based system was limited from the very beginning: The original partition table was installed as part of the boot sector and held space for only four partition entries. These partitions are now called primary partitions.

Logical Partitions

One primary partition of a hard drive may be subpartitioned. These are logical partitions. This effectively allows us to skirt the historical four partition limitation.

The primary partition used to house the logical partitions is called an extended partition and it has its own file system type (0x05). Unlike primary partitions, logical partitions must be contiguous. Each logical partition contains a pointer to the next logical partition, which implies that the number of logical partitions is unlimited. However, linux imposes limits on the total number of any type of partition on a drive, so this effectively limits the number of logical partitions. This is at most 15 partitions total on an SCSI disk and 63 total on an IDE disk.

Swap Partitions

Every process running on your computer is allocated a number of blocks of RAM. These blocks are called pages. The set of in-memory pages which will be referenced by the processor in the very near future is called a "working set." Linux tries to predict these memory accesses (assuming that recently used pages will be used again in the near future) and keeps these pages in RAM if possible.

If you have too many processes running on a machine, the kernel will try to free up RAM by writing pages to disk. This is what swap space is for. It effectively increases the amount of memory you have available. However, disk I/O is about a hundred times slower than reading from and writing to RAM. Consider this emergency memory and not extra memory.

If memory becomes so scarce that the kernel pages out from the working set of one process in order to page in for another, the machine is said to be thrashing. Some readers might have inadvertently experienced this: the hard drive is grinding away like crazy, but the computer is slow to the point of being unusable. Swap space is something you need to have, but it is no substitute for sufficient RAM. See the discussion in the section called "Swap Partitions" for tips on determining the size of swap space you need.

Partitioning requirements

What Partitions do I need?

For the Boot Drive: If you want to boot your operating system from the drive you are about to partition, you will need:

- A primary partition
- One or more swap partitions
- Zero or more primary/logical partitions

For any other drive:

- One or more primary/logical partitions
- Zero or more swap partitions

Discussion:

Boot Partition: Your boot partition ought to be a primary partition, not a logical partition. This will ease recovery in case of disaster, but it is not technically necessary. It must be of type 0x83 "Linux native". If you are using a version of **lilo** [<http://tldp.org/HOWTO/LILO.html>] before 21-3 (ie, from the 1990s), your boot partition must be contained within the first 1024 cylinders of the drive. (Typically, the boot partition need only contain the kernel image.)

If you have more than one boot partition (from other OSs, for example,) keep them all in the first 1024 cylinders (*All* DOS partitions must be within the first 1024). If you are using a modern version of lilo, or a means other than lilo to load your kernel (for example, a boot disk or the **LOADLIN.EXE** MS-DOS based Linux loader), the partition can be anywhere. See the Large-disk [<http://www.tldp.org/HOWTO/Large-Disk-HOWTO.html>] HOWTO for details.

Swap Partition: Unless you swap to files (see the section called "Swap Files") you will need a dedicated swap partition. It must be of type 0x82 "Linux swap". It may be positioned anywhere on the disk (but see the section called "Where should I put my swap space?"). Either a primary or logical partition can be used for swap. More than one swap partition can exist on a drive. 8 total (across drives) are permitted. See notes on swap size below (the section called "Swap Partitions").

Logical Partition: A single primary partition must be used as a container (extended partition) for the logical partitions. The extended partition can go anywhere on the disk. The logical partitions must be contiguous, but needn't fill the extended partition.

File Systems

Which file systems need their own partitions?

Everything in your linux file system can go in the same (single) partition. However, there are circumstances when you may want to restrict the growth of certain file systems. For example, if your mail spool was in the same partition as your root fs and it filled the remaining space in the partition, your computer would basically hang.

/var This fs contains spool directories such as those for mail and printing. In addition, it contains the error log directory. If your machine is a server and develops a chronic error, those msgs can fill the partition. Server computers ought to have /var in a different partition than /.

- `/usr` This is where most executable binaries go. In addition, the kernel source tree goes here, and much documentation.
- `/tmp` Some programs write temporary data files here. Usually, they are quite small. However, if you run computationally intensive jobs, like science or engineering applications, hundreds of megabytes could be required for brief periods of time. In this case, keep `/tmp` in a different partition than `/`.
- `/home` This is where users home directories go. If you do not impose quotas on your users, this ought to be in its own partition.
- `/boot` This is where your kernel images go. See discussion above for placement on old systems.

File lifetimes and backup cycles as partitioning criteria

With ext2, partitioning decisions should be governed by backup considerations and to avoid external fragmentation the section called “Some facts about file systems and fragmentation” from different file lifetimes.

Files have lifetimes. After a file has been created, it will remain some time on the system and then be removed. File lifetime varies greatly throughout the system and is partly dependent on the pathname of the file. For example, files in `/bin`, `/sbin`, `/usr/sbin`, `/usr/bin` and similar directories are likely to have a very long lifetime: many months and above. Files in `/home` are likely to have a medium lifetime: several weeks or so. File in `/var` are usually short lived: Almost no file in `/var/spool/news` will remain longer than a few days, files in `/var/spool/lpd` measure their lifetime in minutes or less.

For backup it is useful if the amount of daily backup is smaller than the capacity of a single backup medium. A daily backup can be a complete backup or an incremental backup.

You can decide to keep your partition sizes small enough that they fit completely onto one backup medium (choose daily full backups). In any case a partition should be small enough that its daily delta (all modified files) fits onto one backup medium (choose incremental backup and expect to change backup media for the weekly/monthly full dump - no unattended operation possible).

Your backup strategy depends on that decision.

When planning and buying disk space, remember to set aside a sufficient amount of money for backup! Unbacked data is worthless! Data reproduction costs are much higher than backup costs for virtually everyone!

For performance it is useful to keep files of different lifetimes on different partitions. This way the short lived files on the news partition may be fragmented very heavily. This has no impact on the performance of the `/` or `/home` partition.

Swap Partitions

How large should my swap space be?

Conventional wisdom creates swap space equal to the amount of RAM.

But keep in mind that this is just a rule of thumb. It is easily possible to create scenarios where programs have extremely large or extremely small working sets (see the section called “Swap Partitions”). For example, a simulation program with a large data set that is accessed in a very random fashion would have almost no noticeable locality of reference in its data segment, so its working set would be quite large.

On the other hand, a graphics program with many simultaneously opened JPEGs, all but one iconified, would have a very large data segment. But image transformations are all done on one single image, most of the memory occupied by the program is not accessed. The same is true for an editor with many editor windows where only one window is being modified at a time. These programs have - if they are designed properly - a very high locality of reference and large parts of them can be kept swapped out without too severe performance impact. A user who never quits programs once launched would want a lot of swap space for the same reason.

Servers typically are configured with more swap space than their desktop counterparts. Even though a given amount of swap is sufficient for its operations, the server might come under transient heavy loads which cause it to page out at a high rate. Some administrators prefer this to the server crashing altogether. In these cases, swap might be several times the size of ram.

How large *can* my swap space be?

Currently, the maximum size of a swap partition is architecture-dependent. For i386, m68k, ARM and PowerPC, it is "officially" 2Gb. It is 128Gb on alpha, 1Gb on sparc, and 3Tb on sparc64. An opteron on the 2.6 kernel can write to a 16 Tb swap partition. For linux kernels 2.1 and earlier, the limit is 128Mb. The partition may be larger than 128 MB, but excess space is never used. If you want more than 128 MB of swap for a 2.1 and earlier kernel, you have to create multiple swap partitions (8 max). After 2.4, 32 swap areas are "officially" possible. See setting up swap for details.

footnote: "official" max swap size: With kernel 2.4, the limit is 64 swap spaces at a maximum of 64Gb each, although this is not reflected in the man page for **mkswap**. With the 64 bit opteron on the 2.6 kernel, 128 swap areas are permitted, each a whopping 16 Tb! (thanks to Peter Chubb for the calculation)

Where should I put my swap space?

The short answer is anywhere is fine. However, if you are interested in extracting as much speed as possible, there are two basic strategies (other than buying more RAM).

- Split the swap space across multiple drives, or at least on the drive you write to least.
- Put each swap partition on the outer tracks.

Here are the considerations:

- If you have a disk with many heads and one with less heads and both are identical in other parameters, the disk with many heads will be faster. Reading data from different heads is fast, since it is purely electronic. Reading data from different tracks is slow, since it involves physically moving the head.

It follows then that writing swap on a separate drive will be faster than moving the head back and forth on a single drive.

- *Placement:* Older disks have the same number of sectors on all tracks. With these disks it will be fastest to put your swap in the middle of the disks, assuming that your disk head will move from a random track towards the swap area.
- Newer disks use ZBR (zone bit recording). They have more sectors on the outer tracks. With a constant number of rpms, this yields a far greater performance on the outer tracks than on the inner ones. Put your swap on the fast tracks. (In general, low-numbered cylinders are associated low partition numbers. However, see Kristian's more recent comments [<http://lissot.net/partition/mapping.html>] on this issue. -Tony)
- *Usage:* Of course your disk head will not move randomly. If you have swap space in the middle of a disk between a constantly busy home partition and an almost unused archive partition, you would be

better of if your swap were near the home partition for even shorter head movements. You would be even better off, if you had your swap on another otherwise unused disk, though.

- *Striping*: Speed can be increased by writing to multiple swap areas simultaneously. Swap spaces with the same priority will be written to like a RAID. See the section called “Multiple Swap Areas”.

Summary: Put your swap on a fast disk with many heads that is not busy doing other things. If you have multiple disks: Split swap and scatter it over all your disks or even different controllers.

Partitioning with fdisk

This section shows you how to actually partition your hard drive with the **fdisk** utility. Linux allows only 4 primary partitions. You can have a much larger number of logical partitions by sub-dividing one of the primary partitions. Only one of the primary partitions can be sub-divided.

Examples:

1. Four primary partitions (see the section called “Four primary partitions”)
2. Mixed primary and logical partitions (see the section called “Mixed primary and logical partitions”)

fdisk usage

fdisk is started by typing (as root) **fdisk device** at the command prompt. *device* might be something like `/dev/hda` or `/dev/sda` (see the section called “Naming Convention”). The basic **fdisk** commands you need are:

p print the partition table

n create a new partition

d delete a partition

q quit without saving changes

w write the new partition table and exit

Changes you make to the partition table do not take effect until you issue the write (w) command. Here is a sample partition table:

```
Disk /dev/hdb: 64 heads, 63 sectors, 621 cylinders
Units = cylinders of 4032 * 512 bytes
```

Device	Boot	Start	End	Blocks	Id	System
/dev/hdb1	*	1	184	370912+	83	Linux
/dev/hdb2		185	368	370944	83	Linux
/dev/hdb3		369	552	370944	83	Linux
/dev/hdb4		553	621	139104	82	Linux swap

The first line shows the geometry of your hard drive. It may not be physically accurate, but you can accept it as though it were. The hard drive in this example is made of 32 double-sided platters with one head on each side (probably not true). Each platter has 621 concentric tracks. A 3-dimensional track (the same track on all disks) is called a cylinder. Each track is divided into 63 sectors. Each sector contains 512 bytes of data. Therefore the block size in the partition table is 64 heads * 63 sectors * 512 bytes er...divided by 1024. (See 4 for discussion on problems with this calculation.) The start and end values are cylinders.

Four primary partitions

The overview:

Decide on the size of your swap space (see the section called “Swap Partitions”) and where it ought to go (see the section called “Where should I put my swap space?”). Divide up the remaining space for the three other partitions.

Example:

I start fdisk from the shell prompt:

```
# fdisk /dev/hdb
```

which indicates that I am using the second drive on my IDE controller. (See the section called “Device names”.) When I print the (empty) partition table, I just get configuration information.

```
Command (m for help): p
```

```
Disk /dev/hdb: 64 heads, 63 sectors, 621 cylinders
Units = cylinders of 4032 * 512 bytes
```

I knew that I had a 1.2Gb drive, but now I really know: $64 * 63 * 512 * 621 = 1281982464$ bytes. I decide to reserve 128Mb of that space for swap, leaving 1153982464. If I use one of my primary partitions for swap, that means I have three left for ext2 partitions. Divided equally, that makes for 384Mb per partition. Now I get to work.

```
Command (m for help): n
```

```
Command action
```

```
  e   extended
```

```
  p   primary partition (1-4)
```

```
p
```

```
Partition number (1-4): 1
```

```
First cylinder (1-621, default 1):<RETURN>
```

```
Using default value 1
```

```
Last cylinder or +size or +sizeM or +sizeK (1-621, default 621): +384M
```

Next, I set up the partition I want to use for swap:

```
Command (m for help): n
```

```
Command action
```

```
  e   extended
```

```
  p   primary partition (1-4)
```

```
p
```

```
Partition number (1-4): 2
```

```
First cylinder (197-621, default 197):<RETURN>
```

```
Using default value 197
```

```
Last cylinder or +size or +sizeM or +sizeK (197-621, default 621): +128M
```

Now the partition table looks like this:

Device	Boot	Start	End	Blocks	Id	System
/dev/hdb1		1	196	395104	83	Linux
/dev/hdb2		197	262	133056	83	Linux

I set up the remaining two partitions the same way I did the first. Finally, I make the first partition bootable:

```
Command (m for help): a
Partition number (1-4): 1
```

And I make the second partition of type swap:

```
Command (m for help): t
Partition number (1-4): 2
Hex code (type L to list codes): 82
Changed system type of partition 2 to 82 (Linux swap)
Command (m for help): p
```

The end result:

Disk /dev/hdb: 64 heads, 63 sectors, 621 cylinders
Units = cylinders of 4032 * 512 bytes

Device	Boot	Start	End	Blocks	Id	System
/dev/hdb1	*	1	196	395104+	83	Linux
/dev/hdb2		197	262	133056	82	Linux swap
/dev/hdb3		263	458	395136	83	Linux
/dev/hdb4		459	621	328608	83	Linux

Finally, I issue the write command (w) to write the table on the disk.

Side topics:

- the section called “Activating Swap Space”
- the section called “Formating Partitions”
- the section called “Mounting Partitions”

Mixed primary and logical partitions

The overview: create one use one of the primary partitions to house all the extra partitions. Then create logical partitions within it. Create the other primary partitions before or after creating the logical partitions.

Example:

I start fdisk from the shell prompt:

```
# fdisk /dev/sda
```

which indicates that I am using the first drive on my SCSI chain. (See the section called “Device names”).

First I figure out how many partitions I want. I know my drive has a 183Gb capacity and I want 26Gb partitions (because I happen to have back-up tapes that are about that size).

183Gb / 26Gb = ~7

so I will need 7 partitions. Even though fdisk accepts partition sizes expressed in Mb and Kb, I decide to calculate the number of cylinders that will end up in each partition because fdisk reports start and stop points in cylinders. I see when I enter fdisk that I have 22800 cylinders.

```
> The number of cylinders for this disk is set to 22800.  There is
> nothing wrong with that, but this is larger than 1024, and could in
> certain setups cause problems with:  1) software that runs at boot
> time (e.g., LILO)  2) booting and partitioning software from other
> OSs  (e.g., DOS FDISK, OS/2 FDISK)
```

So, 22800 total cylinders divided by seven partitions is 3258 cylinders. Each partition will be about 3258 cylinders long. I ignore the warning msg because this is not my boot drive (the section called "Partitioning requirements").

Since I have 4 primary partitions, 3 of them can be 3258 long. The extended partition will have to be (4 * 3258), or 13032, cylinders long in order to contain the 4 logical partitions.

I enter the following commands to set up the first of the 3 primary partitions (stuff I type is bold):

```
Command (m for help): n
Command action
  e   extended
  p   primary partition (1-4)
p
Partition number (1-4): 1
First cylinder (1-22800, default 1): <RETURN>
Using default value 1
Last cylinder or +size or +sizeM or +sizeK (1-22800, default 22800): 3258
```

The last partition is the extended partition:

```
Partition number (1-4): 4
First cylinder (9775-22800, default 9775): <RETURN>
Using default value 9775
Last cylinder or +size or +sizeM or +sizeK (9775-22800, default 22800): <RETURN>
Using default value 22800
```

The result, when I issue the print table command is:

/dev/sda1	1	3258	26169853+	83	Linux
/dev/sda2	3259	6516	26169885	83	Linux
/dev/sda3	6517	9774	26169885	83	Linux
/dev/sda4	9775	22800	104631345	5	Extended

Next I segment the extended partition into 4 logical partitions, starting with the first logical partition, into 3258-cylinder segments. The logical partitions automatically start from /dev/sda5.

```
Command (m for help): n
First cylinder (9775-22800, default 9775): <RETURN>
Using default value 9775
```

Last cylinder or +size or +sizeM or +sizeK (9775-22800, default 22800): 13032

The end result is:

```

Device Boot      Start         End      Blocks   Id  System
/dev/sda1                1         3258    26169853+  83  Linux
/dev/sda2           3259         6516    26169885   83  Linux
/dev/sda3           6517         9774    26169885   83  Linux
/dev/sda4           9775        22800   104631345    5  Extended
/dev/sda5           9775        13032    26169853+  83  Linux
/dev/sda6        13033        16290    26169853+  83  Linux
/dev/sda7        16291        19584    26459023+  83  Linux
/dev/sda8        19585        22800    25832488+  83  Linux

```

Finally, I issue the write command (w) to write the table on the disk. To make the partitions usable, I will have to format (the section called “Formatting Partitions”) each partition and then mount (the section called “Mounting Partitions”) it.

Submitted Examples

I'd like to submit my partition layout, because it works well with any distribution of Linux (even big RPM based ones). I have one hard drive that ... is 10 gigs, exactly. Windows can't see above 9.3 gigs of it, but Linux can see it all, and use it all. It also has much more than 1024 cylinders.

Table 7. Partition layout example

Partition	Mount point	Size
/dev/hda1	/boot	(15 megs)
/dev/hda2	windows 98 partition	(2 gigs)
/dev/hda3	extended	(N/A)
/dev/hda5	swap space	(64 megs)
/dev/hda6	/tmp	(50 megs)
/dev/hda7	/	(150 megs)
/dev/hda8	/usr	(1.5 gigs)
/dev/hda9	/home	(rest of drive)

I test new kernels for the USB mass storage, so that explains the large /boot partition. I install LILO into the MBR, and by default I boot windows (I'm not the only one to use this computer).

I also noticed that you don't have any REAL examples of partition tables, and for newbies I HIGHLY suggest putting quite a few up. I'm freshly out of the newbie stage, and partitioning was what messed me up the most.

Valkor [mailto:valkor@qx.net]

Labels

In linux, hard drives are referred to as devices, and devices are pseudo files in /dev. For example, the first partition of the second lowest numbered SCSI drive is /dev/sdb1. If the drive referred to as /dev/sda is removed from the chain, then the latter partition is automatically renamed /dev/sda1 at reboot.

Volume Labels

Volume labels make it possible for partitions to retain a consistent name regardless of where they are connected, and regardless of whatever else is connected. Labels are not mandatory for a linux volume. Each can be a maximum of 16 characters long.

There are three tools to make volume labels: `mke2fs`, `tune2fs` and `e2label`.

Simple Invocation

```
e2label /dev/hdb1 pubsw
```

```
tune2fs -L pubsw /dev/hdb1
```

Either of these two commands will label the first partition of the second drive "pubsw". That label stays with that particular partition, even if the drive is moved to another controller or even another computer.

```
mke2fs pubsw /dev/hdb1
```

```
mke2fs -L pubsw /dev/hdb1
```

will do the same thing as the first two commands - after they make the file system. This means that either of these last two commands will delete any existing data in the partition.

How to Use

Here is a sample `fstab`. This is a text file located in `/etc`, which is usually set up during the installation of the operating system. It describes where each partition will be mounted, and how it will be mounted. It can be modified by you, either through a utility or manually, when you add/remove devices.

```
LABEL=/          /          ext3    defaults    1 1
LABEL=/boot      /boot      ext2    defaults    1 2
none             /dev/pts   devpts  gid=5,mode=620 0 0
none             /dev/shm   tmpfs   defaults    0 0
LABEL=HOME       /home      ext3    defaults    1 2
none             /proc      proc    defaults    0 0
none             /sys       sysfs   defaults    0 0
LABEL=/usr       /usr       ext3    defaults    1 2
/dev/hdc1        /k-space   ext3    defaults    1 2
/dev/hda6        swap       swap    defaults    0 0
/dev/hdd         /media/cdrecorder auto    pamconsole,ro,exec,noauto,managed 0 0
/dev/fd0         /media/floppy auto    pamconsole,exec,noauto,managed 0 0
```

The leftmost column lists devices and the second column lists mount points. This example contains a mixture of devices and labels. The master drive of the second controller is always mounted on `/k-space`. The partition labeled "HOME" is always mounted on `/home`, regardless of which drive it is on or which partition number it has. Notice that it is permissible to use mount points as labels, such as `"/usr"`

Device Labels

`devlabel` is a script which creates symbolic links to devices. For example,

```
devlabel -d /dev/hdb1 -s /dev/home
```


will create a link from `/dev/hdb1` to `/dev/home`. Crucially, it stores a unique identifier for the hardware that was on `/dev/hdb1` and stores that identifier along with the link name that you specified in `/etc/sysconfig/devlabel`. If the hardware is later moved to `/dev/hdc1`, its unique identifier will be queried (using `/usr/bin/partition_uuid`), matched to its entry in `/etc/sysconfig/devlabel`, and again linked to `/dev/home`.

Formatting an ext2/3 partition

When a hard drive is partitioned, it is mapped into sections, but the sections are empty. It is like a newly constructed library; shelves, signs, and a card catalogue system must be put in place before the books are put away.

The organizational structure inside a partition is called a file system. With Linux, the standard file system is ext2 and ext3. The ext3 file system is ext2, plus a log of disk writes called a journal. The journal allows the system to recover quickly from accidental power outages, among other things.

The principal tool for making an ext2/3 file system in a partition is **mke2fs**. It is usually found in `/sbin`. **mkfs.ext2** and **mkfs.ext3** are frontends which pass specific options to **mke2fs**.

Simple Invocation

```
mke2fs /dev/hdb1
```

```
mkfs.ext2 /dev/hdb1
```

both of which make an ext2 file system on the first partition of the second drive, and

```
mke2fs -j /dev/hdb1
```

```
mkfs.ext3 /dev/hdb1
```

make an ext3 file system.

Reserved blocks

The `-m` option is probably the one of most use to non-experts. If the file system becomes filled and there is no more space to write, it is basically unusable because the operating system is constantly writing to disk. By default, five percent of the partition is reserved for use by the root user. This allows root to conduct administrative activities on the partition and perhaps move some data off. However, this is most critical when the partition contains `/` or `home` directories. For pure data partitions, this is just lost space. Five percent of a 250Gb partition is 12.5 Gb. Especially in the case of large partitions, it is safe to set the reserved space to the minimum, which is one percent.

```
mkfs.ext3 -m 1 /dev/hdb1
```

creates a file system with only 1% of its space reserved for the root user. `tune2fs -m` can be used to adjust the reserved blocks after data is loaded on the partition.

Recovering a Deleted Partition Table

Below are instructions for manually recovering a deleted partition table. There are utilities such as `gpart` [<http://www.stud.uni-hannover.de/user/76201/gpart>] or `TestDisk` [<http://www.cgsecurity.org/index.html?testdisk.html>] which can make this task considerably easier. If you are reading this, however, because you have run out of luck, this is what you will have to do:

1. Make a partition that is at least as big as your first partition was. You can make it larger than the original partition by any amount. If you underestimate, there will be much wailing and gnashing of teeth.

```

Command (m for help): n
Command action
  e   extended
  p   primary partition (1-4)
P
Partition number (1-4): 1
First cylinder (1-23361, default 1): <RETURN>
Using default value 1
Last cylinder or +size or +sizeM or +sizeK (1-22800, default 22800): 13032

Command (m for help): w

```

2. Run **dumpe2fs** on the first partition and grep out the block count.

Example:

```

% dumpe2fs /dev/sda1 | grep "Block count:"
Block count:                41270953

```

If you are uncertain about this value, repeat Step 1 with a bigger partition size. If the block count changes, then you underestimated the size of the original partition. Repeat Step 1 until you get a stable block count.

3. Remove the partition you just created

```

Command (m for help): d
Partition number (1-4): 1

```

4. Make a new partition with the exact size you got from the block count. Since you cannot enter block size in fdisk, you need to figure out how many cylinders to request. Here is the formula:

$$(\text{number of needed cylinders}) = (\text{number of blocks}) / (\text{block size})$$

$$(\text{block size}) = (\text{unit size}) / 1024$$

$$(\text{unit size}) = (\text{number of heads}) * (\text{number of sectors/cylinder}) * (\text{number of bytes/sector})$$

Consider the following example, where a hard drive has been partitioned into four primary partitions of 1, 2, 4, and 8 cylinders.

```

disk /dev/sda: 16 heads, 63 sectors, 23361 cylinders
Units = cylinders of 1008 * 512 bytes

```

Device	Boot	Start	End	Blocks	Id	System
/dev/sda1		1	2	976+	83	Linux
/dev/sda2		3	5	1512	83	Linux

```
/dev/sda3          6          10          2520    83  Linux
/dev/sda4          11         19          4536    83  Linux
```

fdisk provides the configuration information I need in the head of the output. The unit size is **516096** (**16** heads * **63** sectors/cyl * **512** bytes/sector). The block size is **504** (**516096** / **1024**). The number of needed cylinders for the second partition is therefore **3** (**1512** blocks / **504**). The partition table shows that this is indeed the case: the first cylinder is **3**, the second **4**, and the last is **5**, for a total of three cylinders. The number of needed cylinders for the third partition is calculated similarly: **2520** blocks / **504** = **5**, which corresponds to blocks **6,7,8,9,10** . Notice that this calculation does not work for the first partition because the block count is wrong (**976** instead of **1008**). The plus sign indicates that not all the blocks are included in the **fdisk** value. When you try the calculation (**976** / **504**) you get **1.937**. Knowing that the number of cylinders must be an integer, you can simply round up.

5. Run **e2fsck** on it to verify that you can read the new partition.

6. Repeat Steps 1-5 on remaining partitions.

Remount your partitions. Amazingly, all of your data will be there.

Credit goes to: Mike Vevea, jedi sys admin, for providing the basic strategy.

Setting Up Swap Space

Swap Files

Normally, there are only two steps to setting up swap space, creating the partition and adding it to `/etc/fstab`. A typical `fstab` entry for a swap partition at `/dev/hda6` would look like this:

```
/dev/hda6 swap swap defaults 0 0
```

The next time you reboot, the initialization scripts will activate it automatically and there's nothing more to be done.

However, if you want to make use of it right away, you'll need to activate it manually. As root, type:

```
mkswap -f /dev/hda6
swapon /dev/hda6
```

Swap Files

There might be times when you've run out of swap space and it is not practical to repartition a drive or add a new one. In this case, you can use a regular file in an ordinary partition. All you have to do is create a file of the size you want

```
dd if=/dev/zero of=/var/my_swap bs=1024 count=131072
```

and activate it

```
mkswap -f /var/my_swap
swapon /var/my_swap
```

This invocation creates a file called `my_swap` in `/var`. It is 128 Mb long ($128 \times 1024 = 131072$). Initially, it is filled with zeros. However, **mkswap** marks it as swap space and **swapon** tells the kernel to start using it as swap space. When you are done with it,

```
swapoff /var/my_swap
rm /var/my_swap
```

Multiple Swap Areas

More than one swap partition can be used on the same system. Consider an example `fstab` where there is a single swap partition:

```
/dev/hda5  /          ext3  defaults        1 1
/dev/hda1  /boot      ext2  defaults        1 2
none      /dev/pts  devpts gid=5,mode=620  0 0
none      /proc      proc  defaults        0 0
/dev/hda7  /usr       ext3  defaults        1 2
/dev/hda6  swap       swap  defaults        0 0
```

Imagine replacing the entry for the swap partition with these three lines:

```
/dev/hda6  none      swap    sw,pri=3       0 0
/dev/hdb2  none      swap    sw,pri=2       0 0
/dev/hdc2  none      swap    sw,pri=1       0 0
```

This configuration would cause the kernel to use `/dev/hda6` first. it has the highest priority assigned to it (`pri=3`). The maximum priority can be 32767 and the lowest 0. If that space were to max out, the kernel would start using `/dev/hdb2`, and on to `/dev/hdc2` after that. Why such a configuration? Imagine that the newest (fastest) drives are given the highest priority. This will minimize speed loss as swap space usage grows.

It is possible to write to all three simultaneously. If each has the same priority, the kernel will write to them much like a RAID, with commensurate speed increases.

```
/dev/hda6  none      swap    sw,pri=3       0 0
/dev/hdb2  none      swap    sw,pri=3       0 0
/dev/hdc2  none      swap    sw,pri=3       0 0
```

Notice that these three partitions are on separate drives, which is ideal in terms of speed enhancement.

Appendix

Formating Partitions

At the shell prompt, I begin making the file systems on my partitions. Continuing with the example in (see the section called “Mixed primary and logical partitions”), this is:

```
# mke2fs /dev/sda1
```

I need to do this for each of my partitions, but not for `/dev/sda4` (my extended partition). Linux supports types of file systems other than `ext2`. You can find out what kinds your kernel supports by looking in: `/usr/src/linux/include/linux/fs.h`

The most common file systems can be made with programs in `/sbin` that start with "mk" like **mkfs.msdos** and **mke2fs**.

Activating Swap Space

To set up a swap partition:

```
# mkswap -f /dev/hda5
```

To activate the swap area:

```
# swapon /dev/hda5
```

Normally, the swap area is activated by the initialization scripts at boot time.

Mounting Partitions

Mounting a partition means attaching it to the linux file system. To mount a linux partition:

```
# mount -t ext2 /dev/sda1 /opt
```

`-t ext2` File system type. Other types you are likely to use are:

- `ext3` (journaling file system based on `ext2`)
- `msdos` (DOS)
- `hfs` (mac)
- `iso9660` (CDROM)
- `nfs` (network file system)

`/dev/sda1` Device name. Other device names you are likely to use:

- `/dev/hdb2` (second partition in second IDE drive)
- `/dev/fd0` (floppy drive A)
- `/dev/cdrom` (CDROM)

`/opt` mount point. This is where you want to "see" your partition. When you type `ls /opt`, you can see what is in `/dev/sda1`. If there are already some directories and/or files under `/opt`, they will be invisible after this mount command.

Some facts about file systems and fragmentation

Disk space is administered by the operating system in units of blocks and fragments of blocks. In `ext2`, fragments and blocks have to be of the same size, so we can limit our discussion to blocks.

Files come in any size. They don't end on block boundaries. So with every file a part of the last block of every file is wasted. Assuming that file sizes are random, there is approximately a half block of waste for each file on your disk. Tanenbaum calls this "internal fragmentation" in his book "Operating Systems".

You can guess the number of files on your disk by the number of allocated inodes on a disk. On my disk

```
# df -i
Filesystem          Inodes    IUsed    IFree   %IUsed Mounted on
/dev/hda3           64256    12234    52022    19%   /
/dev/hda5           96000    43058    52942    45%   /var
```

there are about 12000 files on / and about 44000 files on /var. At a block size of 1 KB, about $6+22 = 28$ MB of disk space are lost in the tail blocks of files. Had I chosen a block size of 4 KB, I had lost 4 times this space.

Data transfer is faster for large contiguous chunks of data, though. That's why ext2 tries to preallocate space in units of 8 contiguous blocks for growing files. Unused preallocation is released when the file is closed, so no space is wasted.

Noncontiguous placement of blocks in a file is bad for performance, since files are often accessed in a sequential manner. It forces the operating system to split a disk access and the disk to move the head. This is called "external fragmentation" or simply "fragmentation" and is a common problem with MS-DOS file systems. In conjunction with the abysmal buffer cache used by MS-DOS, the effects of file fragmentation on performance are very noticeable. DOS users are accustomed to defragging their disks every few weeks and some have even developed some ritualistic beliefs regarding defragmentation.

None of these habits should be carried over to Linux and ext2. Linux native file systems do not need defragmentation under normal use and this includes any condition with at least 5% of free space on a disk. There is a defragmentation tool for ext2 called defrag, but users are cautioned against casual use. A power outage during such an operation can trash your file system. Since you need to back up your data anyway, simply writing back from your copy will do the job.

The MS-DOS file system is also known to lose large amounts of disk space due to internal fragmentation. For partitions larger than 256 MB, DOS block sizes grow so large that they are no longer useful (This has been corrected to some extent with FAT32). Ext2 does not force you to choose large blocks for large file systems, except for very large file systems in the 0.5 TB range (that's terabytes with 1 TB equaling 1024 GB) and above, where small block sizes become inefficient. So unlike DOS there is no need to split up large disks into multiple partitions to keep block size down.

Use a 1Kb block size if you have many small files. For large partitions, 4Kb blocks are fine.