



Introduction to Linux for Bioinformatics

Managing data

Joachim Jacob
5 and 12 May 2014



This presentation is available under the Creative Commons Attribution-ShareAlike 3.0 Unported License. Please refer to <http://www.bits.vib.be/> if you use this presentation or parts hereof.



Bioinformatics data

Historically, bioinformatics has always used text files to store data.

Homo sapiens CD99 molecule, mRNA (cDNA clone MGC:19734 IMAGE:3606974), complete cds

GenBank: BC010109.2

[FASTA](#) [Graphics](#)

[Go to](#) ☒

LOCUS BC010109 1255 bp mRNA linear PRI 24-JUL-2006
DEFINITION Homo sapiens CD99 molecule, mRNA (cDNA clone MGC:19734
IMAGE:3606974), complete cds.

ACCESSION BC010109
VERSION BC010109.2 GI:33991438

KEYWORDS MGC.
SOURCE Homo sapiens (human)
ORGANISM [Homo sapiens](#)

Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;
Catarrhini; Hominidae; Homo.
REFERENCE 1 (bases 1 to 1255)

AUTHORS Strausberg, P.
Klausner, R.
Altschul, S.
Hopkins, R.
Diatchenko, F.
Stapleton, M.
Scheetz, T.
Carninci, P.
Abramson, R.
McKernan, K.

HMMS3/b [3.0 | March 2010]

NAME Histones model

LENG 251

ALPH amino

RF np

CS no

MAP yes

DATE Tue Feb 28 14:11:20 2012

NSIQ 6

EFFN 1.004883

CKSUM 427625644

STATS LOCAL MSV -11.2973 0.70287

STATS LOCAL VITERBI -12.0628 0.70287

STATS LOCAL FORWARD -4.6693 0.70287

HMM

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

E

F

G

H

I

K

S

T

V

W

Y

NGS data

The NGS machines produce a lot of data, stored in **plain text** files. These files are multiple gigabytes in size.



Tips for managing (NGS) data

1. *When you move the data, do it in its smallest form.*

→ **Compress** the data.

2. *When you decompress the data, leave it where it is.*

→ **Symbolic links** point to the data in different folders.

3. *Provide enough storage for your data.*

→ choose your **file system type** wisely

Compression: tools in Linux

Compression Tools

Lrzip	Achieve very high compression ratios and speed when used with large files
lbzip2	Multi-threaded implementation of bzip2, suited for serial and parallel processing
7-Zip	File archiver with a high compression ratio
XZ Utils	Successor to the Lempel-Ziv/Markov-chain Algorithm compression format
bzip2	Freely available, patent free, high-quality data compressor
gzip	Provides the standard GNU file compression utilities
PeaZip	Cross-platform portable file archiver

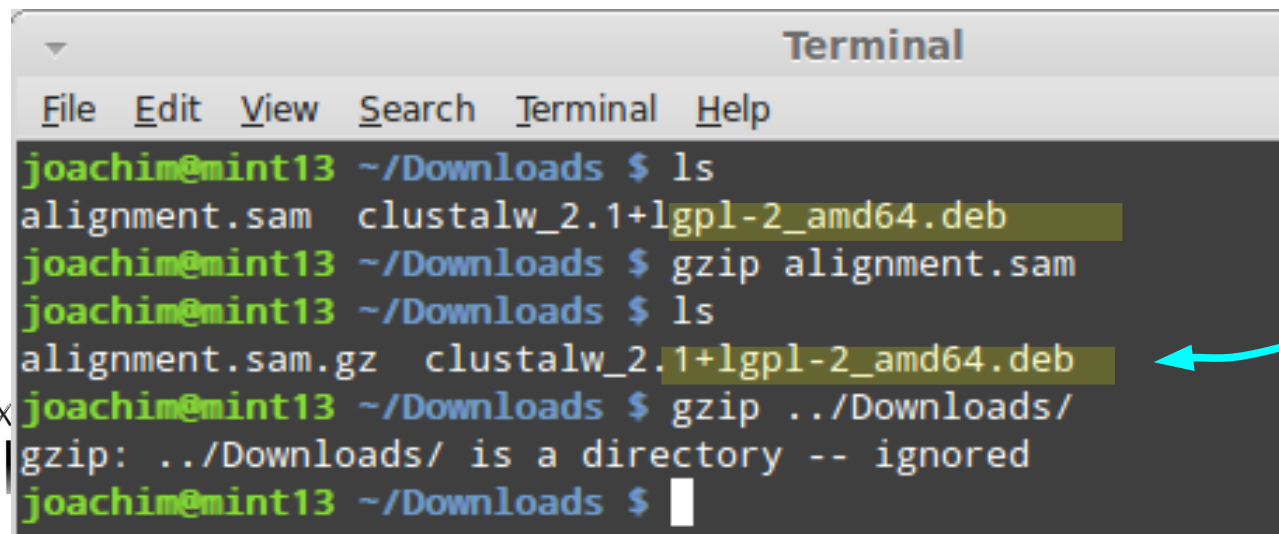
But some more exist (e.g. freearc, dtrx). **gzip** and **bzip2** are the most used and fairly performant.

Tips

Widely used compression tools:

- GNU zip (`gzip`)
- Block Sorting compression (`bzip2`)

Typically, compression tools work on one file.
How to compress complete directories and their contents?



```
Terminal
File Edit View Search Terminal Help
joachim@mint13 ~/Downloads $ ls
alignment.sam  clustalw_2.1+lgpl-2_amd64.deb
joachim@mint13 ~/Downloads $ gzip alignment.sam
joachim@mint13 ~/Downloads $ ls
alignment.sam.gz  clustalw_2.1+lgpl-2_amd64.deb
joachim@mint13 ~/Downloads $ gzip ../Downloads/
gzip: ../Downloads/ is a directory -- ignored
joachim@mint13 ~/Downloads $
```

A terminal window titled "Terminal" with a menu bar (File, Edit, View, Search, Terminal, Help). The user "joachim" is at host "mint13" in the directory "~/Downloads". They run `ls` showing `alignment.sam` and `clustalw_2.1+lgpl-2_amd64.deb`. Then they run `gzip alignment.sam`. They run `ls` again, showing `alignment.sam.gz` and `clustalw_2.1+lgpl-2_amd64.deb`. Finally, they run `gzip ../Downloads/`, which results in the error message `gzip: ../Downloads/ is a directory -- ignored`. A blue arrow points from the right side of the slide to the error message.

Tar without compression



Tar (Tape Archive) is a tool for **bundling a set of files or directories into a single archive**. The resulting file is called a tar ball.

Syntax to create a tarball:

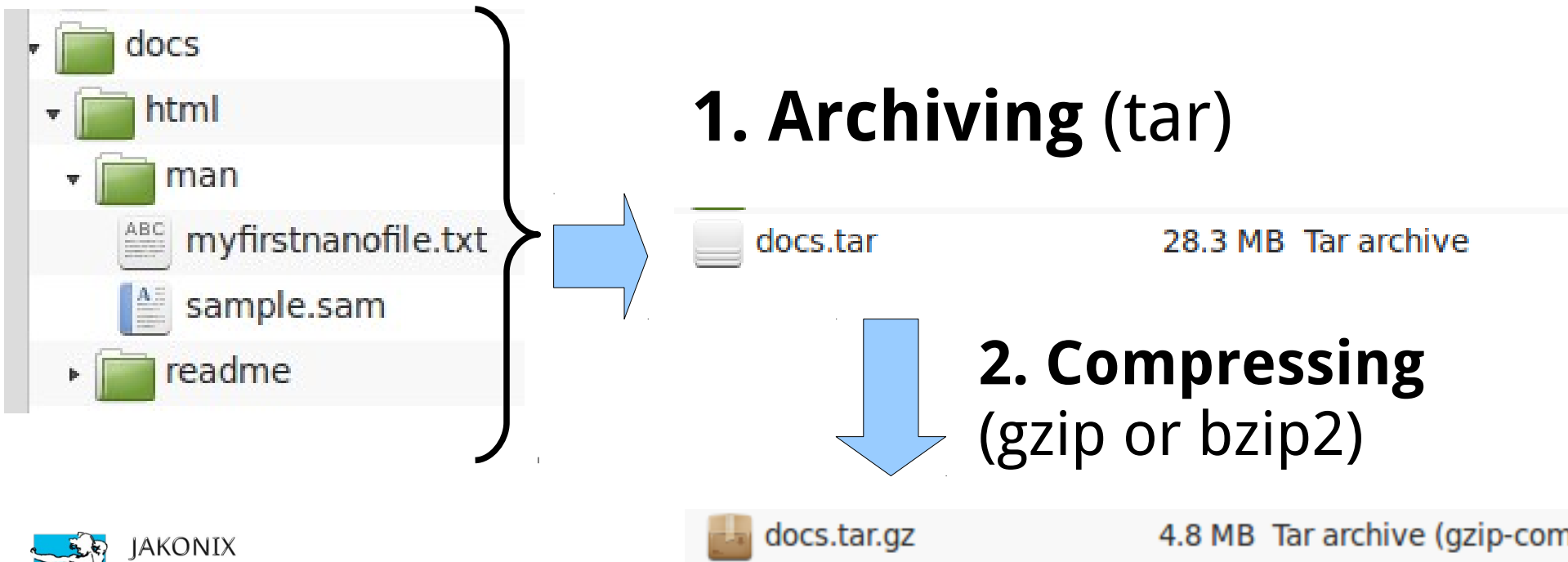
```
$ tar -cf archive.tar file1 file2
```

Syntax to extract:

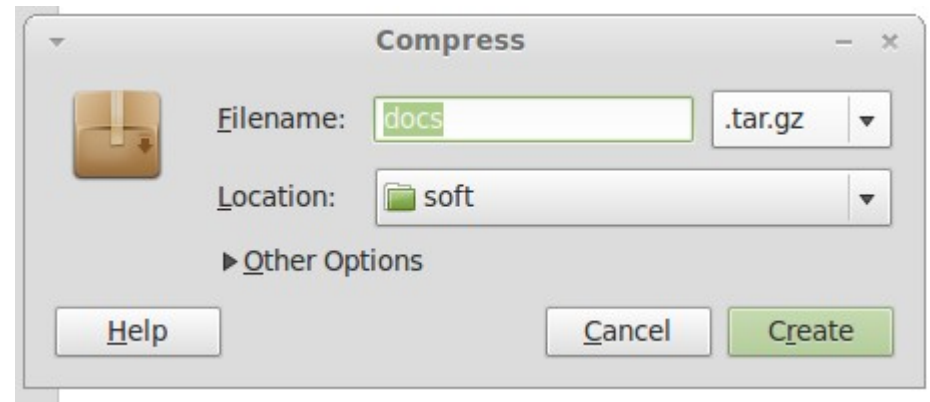
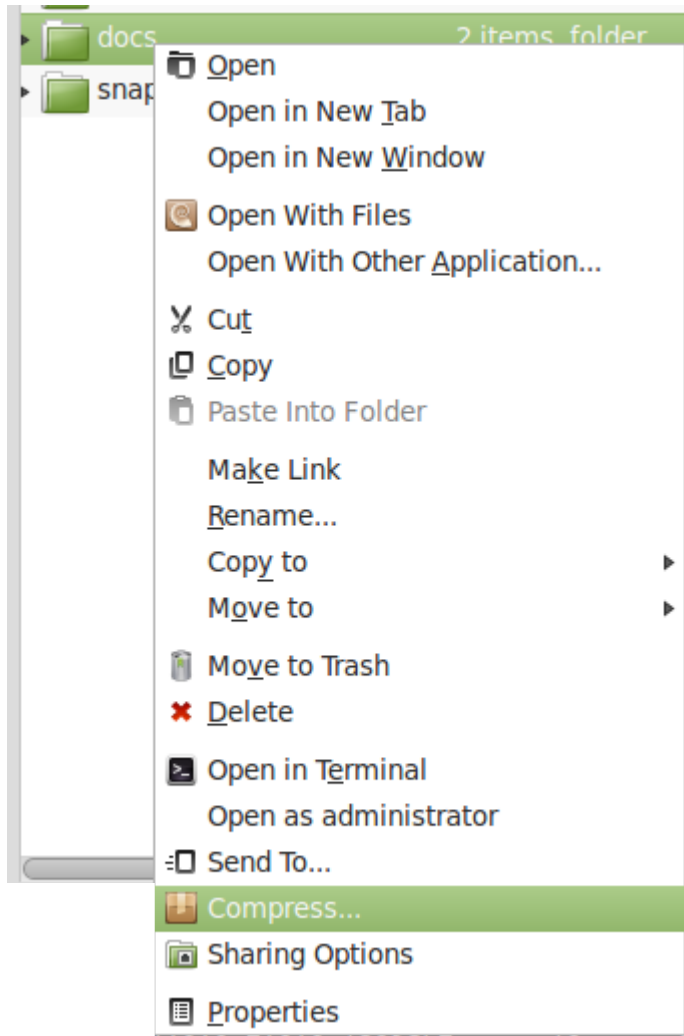
```
$ tar -xvf /path/to/archive.tar
```

Compression: a typical case

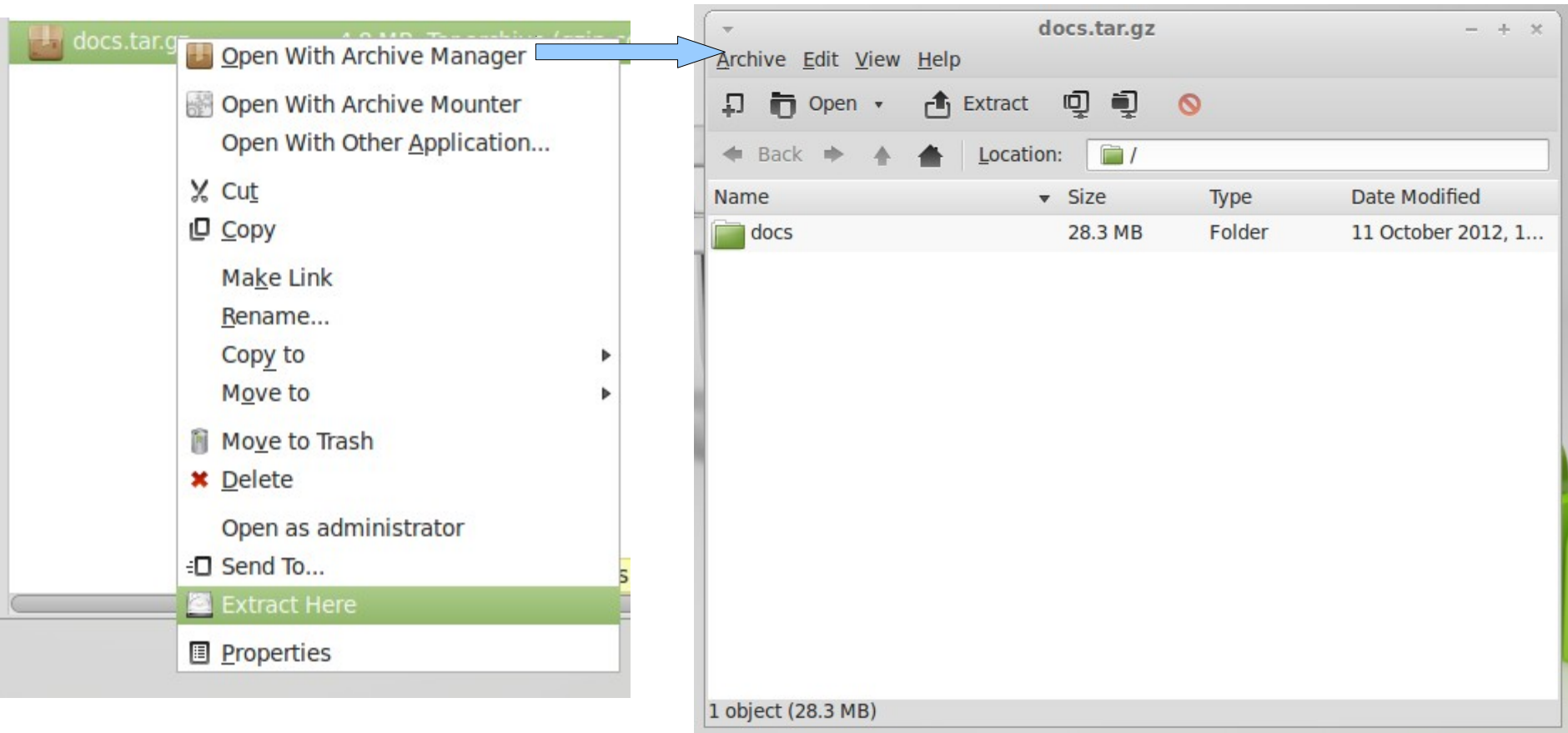
Archiving and compression mostly occur together. The most used formats are **tar.gz** or **tar.bz**. These files are the result of **two** processes.



Compression: on your desktop



Compression: on your desktop



Compression: on the command line

Tar is the tool for creating .tar archives, but it can compress in one go, with the z or j option.

Creating a compressed tar archive:

```
$ tar cvfz mytararchive.tar.gz docs/  
$ tar cvfj mytararchive.tar.bz docs/
```

create

Compression technique

Decompressing a compressed tar archive

```
$ tar xvfz mytararchive.tar.gz  
$ tar xvfj mytararchive.tar.bz
```

extract

files

verbose

De-/compression

To compress one or more files:

```
$ gzip [options] file
```

```
$ bzip2 [options] file
```

To decompress one or more files:

```
$ gunzip [options] file(s)
```

```
$ bunzip2 [options] file(s)
```

Every file will be compressed, and tar.gz or tar.bz appended to it.

Tips

1. Do you have to uncompress a big text file to read it? No! Some tools allow to **read compressed files** (instead of first unpacking then reading). Time saver!

```
$ zcat file(s)  
$ bzip2 file(s)
```

2. Compression is always a **balance** between time and compression ratio. Gzip is faster, bzip2 compresses harder.

If compression is important to you: benchmark!

Exercise

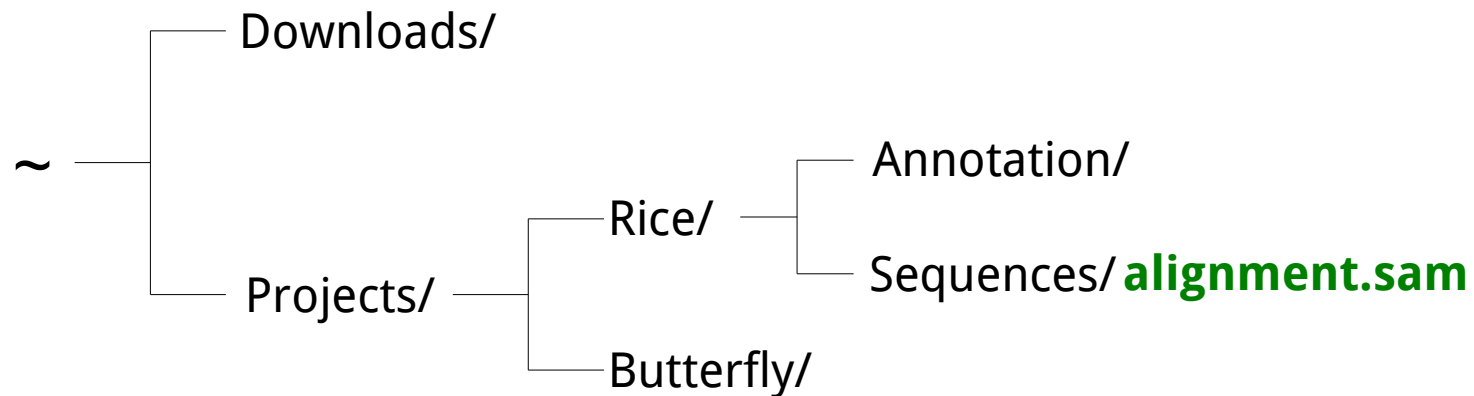


→ a little compression exercise.

Symlinks

Something very convenient!

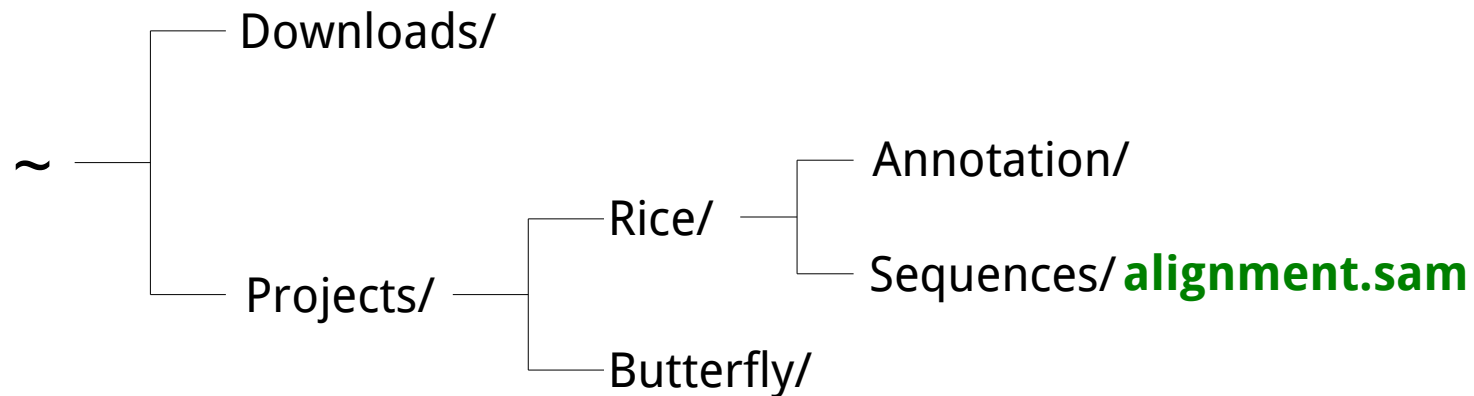
A **symbolic link** (or symlink) is a *file* which points to the location of another file. You can do anything with the symlink that *you can do on the original file*. But when you move the original file from its location, the symlink is 'dead'.



Symlinks

To create a symlink, move to the folder in where the symlink must be created, and execute `ln`.

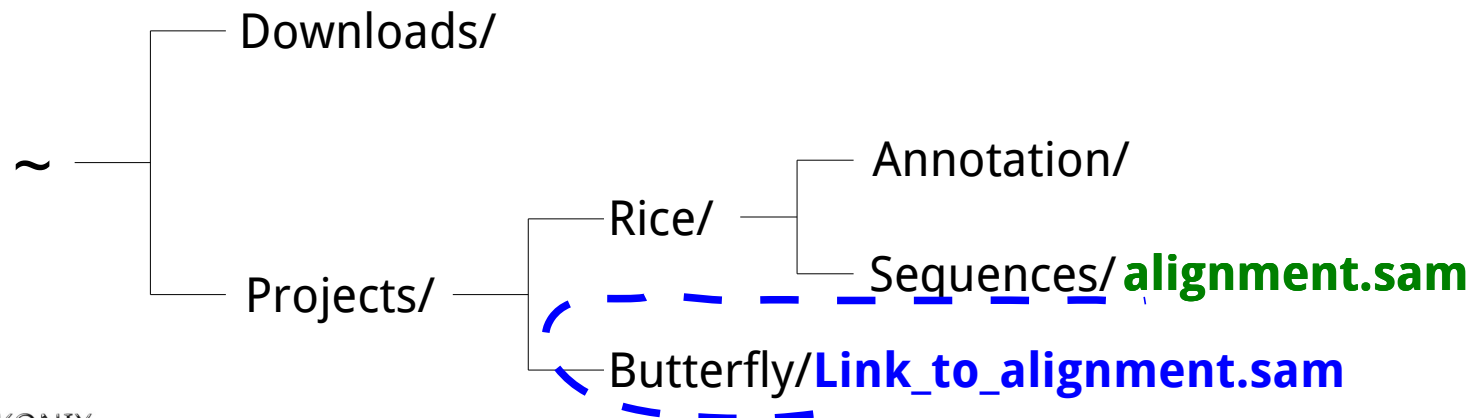
```
~/Projects $ cd Butterfly  
~/Butterfly $ ln -s ../Rice/Sequences/alignment.sam  
Link_to_alignment.sam
```



Symlinks

The symlink is created. You can check with `ls`.
To delete a symlink, use `unlink`.

```
~/Projects $ cd Butterfly
~/Butterfly $ ln -s ../Rice/Sequences/alignment.sam
Link_to_alignment.sam
~/Butterfly $ ls -lh Link_to_alignment.sam
lrwxrwxrwx 1 joachim joachim 44 Oct 22 14:47
Link_to_alignment.sam -> ../Sequences/alignment.sam
```



Exercise



→ a little symlink exercise

Disks and storage

If you dive into bioinformatics, you will have to manage disks and storage.

Two types of disks

- **solid state disks**

Low capacity, high speed, random writes



- **spinning hard disks**

High capacity, 'normal' speed, sequential writes.



A disk is a device

Via the terminal, show the disks using

```
$ sudo fdisk -l
```

```
[sudo] password for joachim:
```

```
Disk /dev/sda: 13.4 GB, 13408141312  
bytes
```

```
...
```

```
Disk /dev/sdb: 3997 MB, 3997171712 bytes
```

```
...
```

A disk is divided into partitions

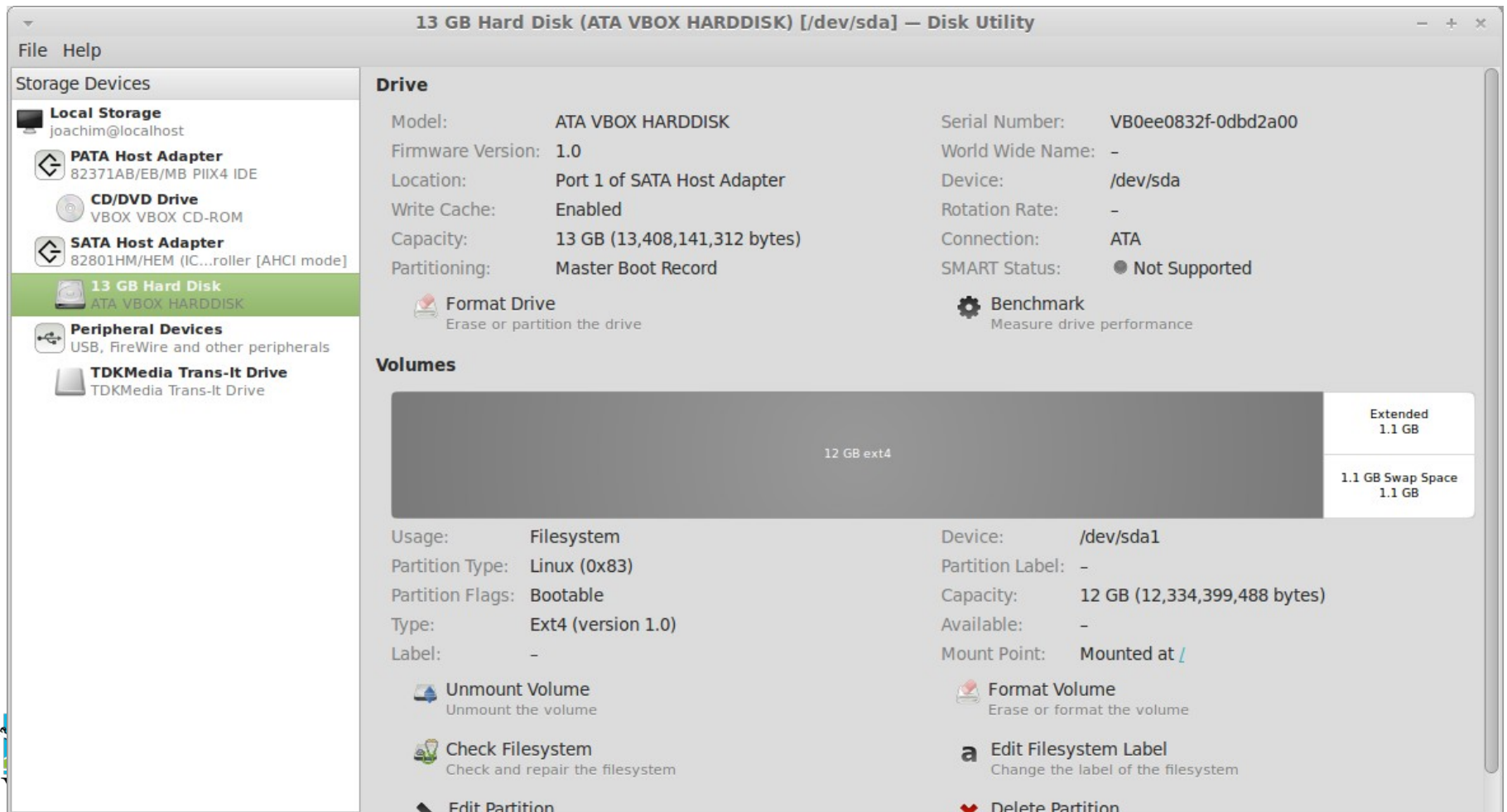
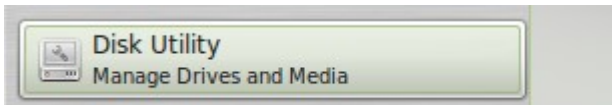


A disk can be divided in parts, called partitions.

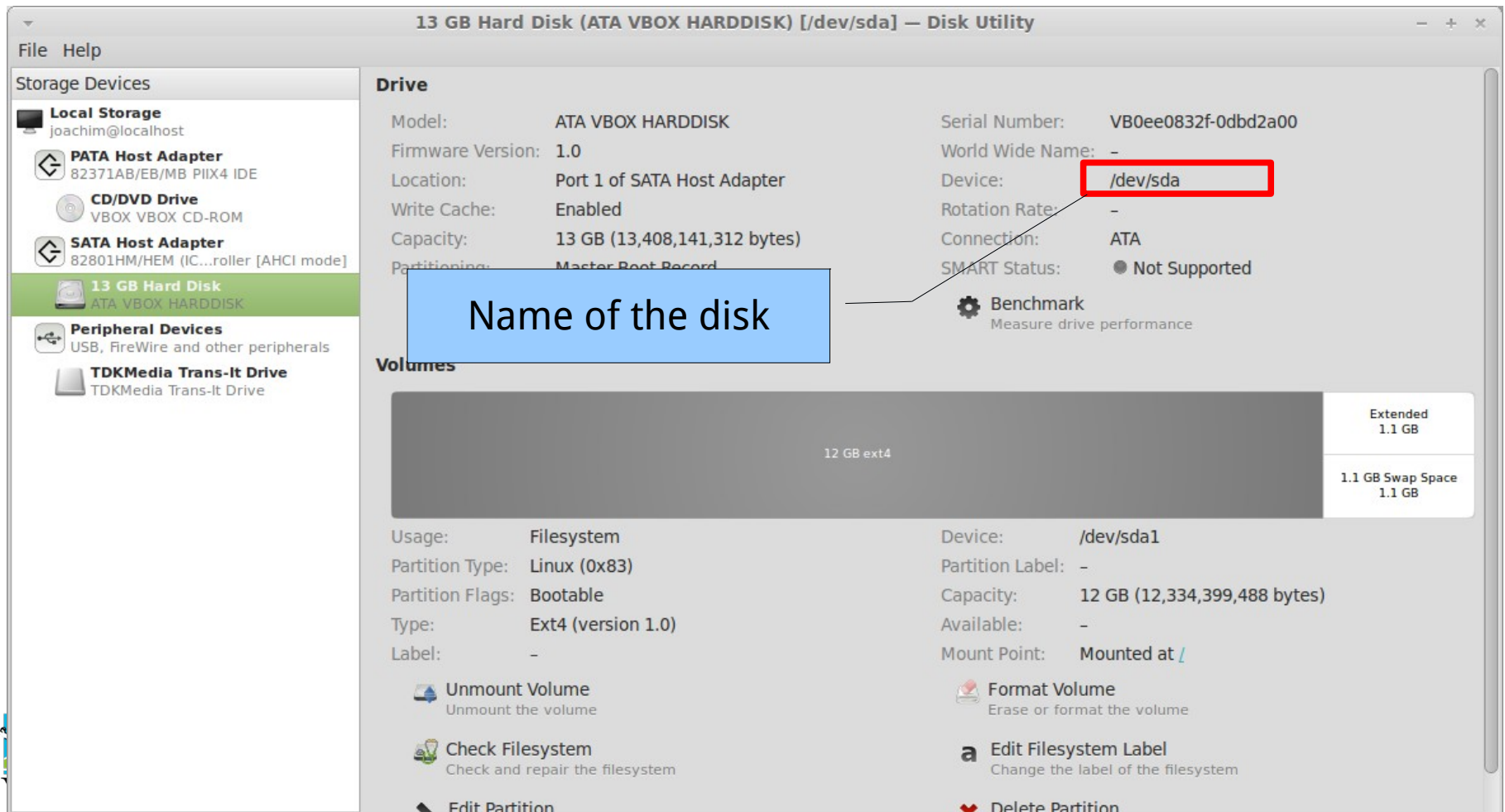
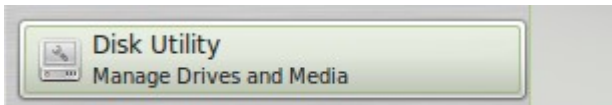
An **internal disk** which runs an operating system is usually divided in partitions, one for each functions.

An **external disk** is usually not divided in partitions.
(but it can be partioned).

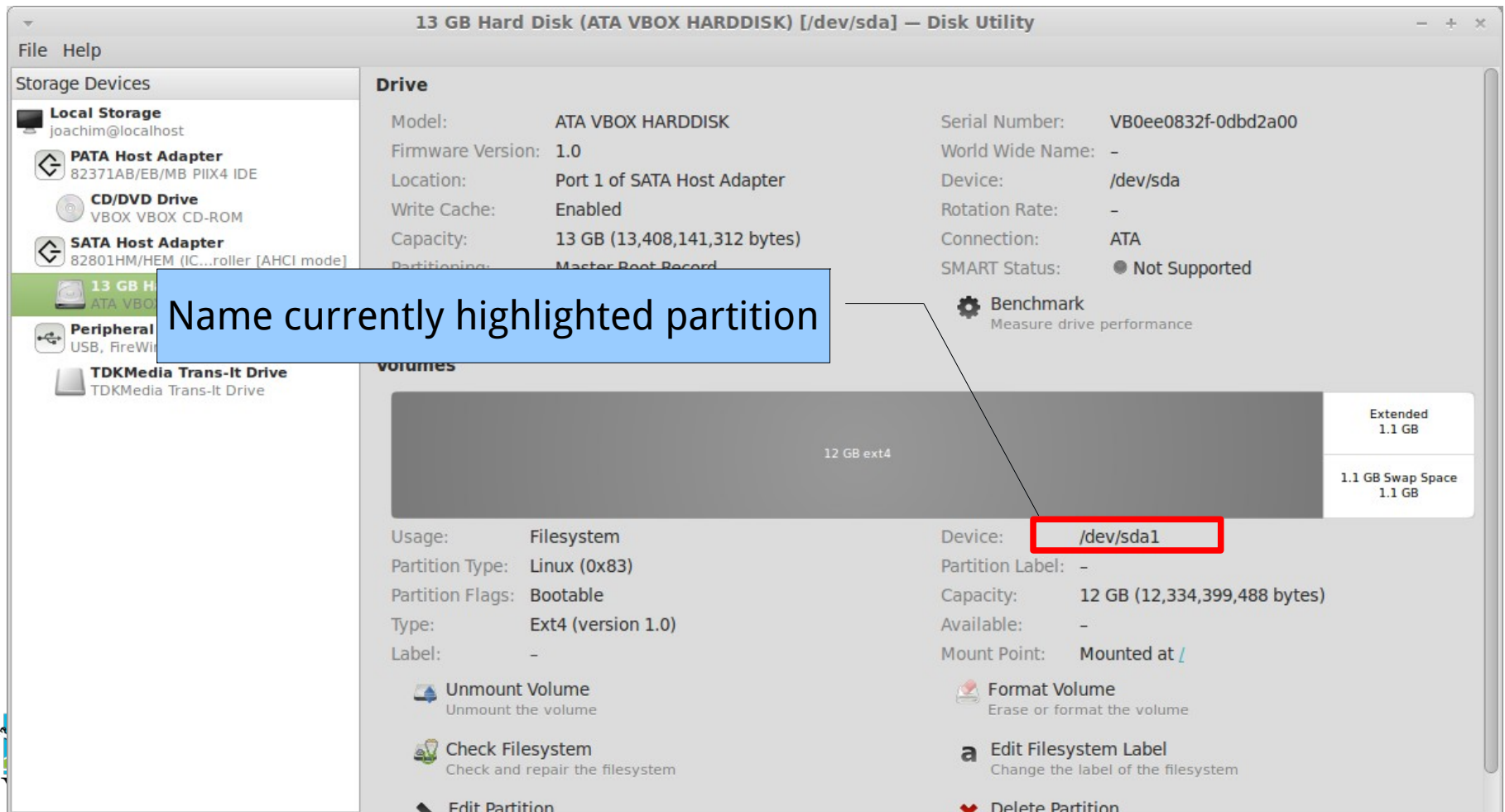
Check out the disk utility tool



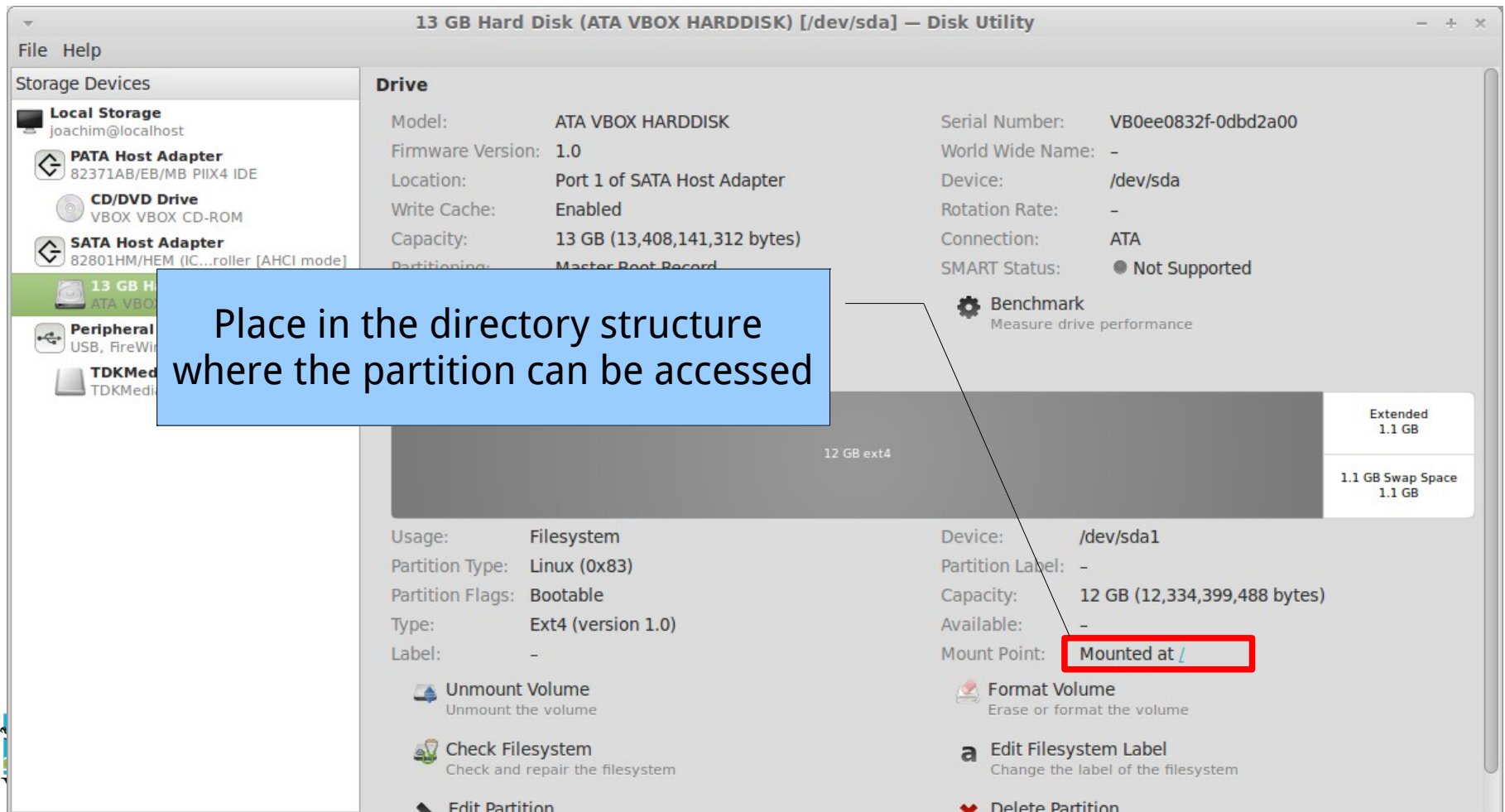
The system disk



The system disk



The system disk



An example of an USB disk

The screenshot shows the 'Disk Utility' window for a 'TDKMedia Trans-It Drive (/dev/sdb)'. The left sidebar lists 'Storage Devices' including Local Storage, PATA Host Adapter, CD/DVD Drive, SATA Host Adapter, 13 GB Hard Disk, and Peripheral Devices. The main area displays drive information: Model (TDKMedia Trans-It Drive), Firmware Version (PMAP), Location (-), Write Cache (-), Capacity (4.0 GB), and Partitioning (Master Boot Record). A 'Format Drive' button is present. On the right, serial number, world wide name, device path (/dev/sdb), rotation rate, connection (USB at 480.0 Mb/s), and SMART status (Not Supported) are shown. A 'Safe Removal' button is also visible. Below the drive information, a partition 'linuxlive' (4.0 GB FAT) is shown. Its details include Usage (Filesystem), Partition Type (W95 FAT32 (LBA) (0x0c)), Partition Flags (Bootable), Type (FAT (32-bit version)), and Label (linuxlive). The 'Mount Point' is highlighted with a red box and labeled 'Mounted at /media/linuxlive'. A blue callout box points to this mount point with the text: 'Place in the directory structure where the partition can be accessed'.

TDKMedia Trans-It Drive (TDKMedia Trans-It Drive) [/dev/sdb] — Disk Utility

Storage Devices

- Local Storage
joachim@localhost
- PATA Host Adapter
82371AB/EB/MB PIIX4 IDE
- CD/DVD Drive
VBOX VBOX CD-ROM
- SATA Host Adapter
82801HM/HEM (IC...roller [AHCI mode])
- 13 GB Hard Disk
ATA VBOX HARDDISK
- Peripheral Devices
USB, FireWire and other peripherals
- TDKMedia Trans-It Drive
TDKMedia Trans-It Drive

Drive

Model: TDKMedia Trans-It Drive
Firmware Version: PMAP
Location: -
Write Cache: -
Capacity: 4.0 GB (3,997,171,712 bytes)
Partitioning: Master Boot Record

Serial Number: 07A1070890FFF8DC
World Wide Name: -
Device: /dev/sdb
Rotation Rate: -
Connection: USB at 480.0 Mb/s
SMART Status: ● Not Supported

Format Drive
Erase or partition the drive

Safe Removal
Power down the drive so it can be removed

linuxlive
4.0 GB FAT

Usage: Filesystem
Partition Type: W95 FAT32 (LBA) (0x0c)
Partition Flags: Bootable
Type: FAT (32-bit version)
Label: linuxlive

Device: /dev/sdb1
Partition Label: -
Capacity: 4.0 GB (3,993,042,944 bytes)
Available: -
Mount Point: Mounted at [/media/linuxlive](#)

Unmount Volume
Unmount the volume

Check Filesystem
Check and repair the filesystem

Delete Partition
Delete the partition

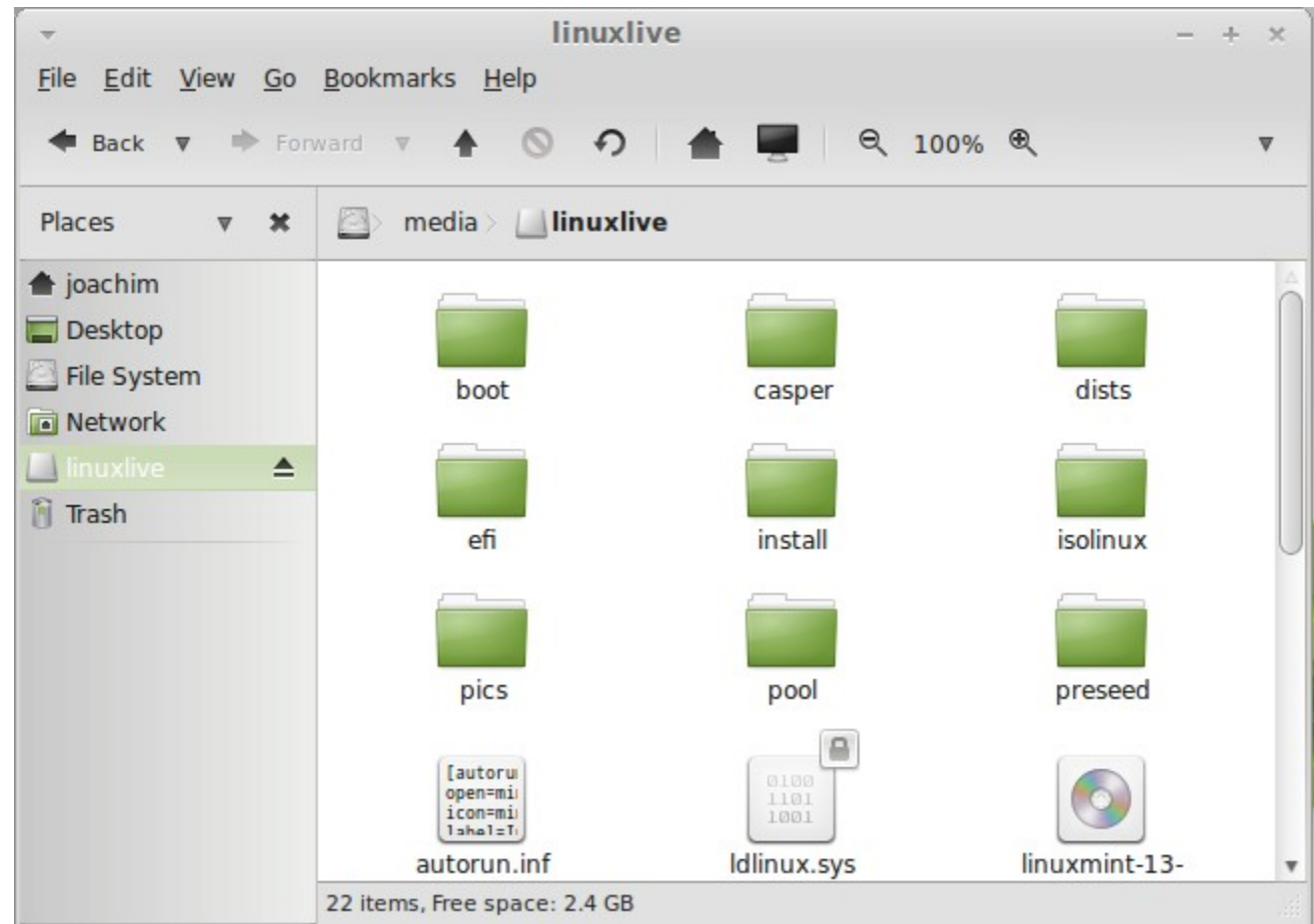
Format Volume
Erase or format the volume

Edit Partition
Change partition type, label and flags

Place in the directory structure where the partition can be accessed

An example of an USB disk

The USB disk is 'mounted' automatically on the directory tree under **/media**.



An example of an USB disk

This is the type of **file system** on the partition (i.e. the way data is stored on this partition)

The partition is said to be **formatted** in FAT32 (in this case).

The screenshot shows the 'Disk Utility' window for a TDKMedia Trans-It Drive. The window title is 'TDKMedia Trans-It Drive (TDKMedia Trans-It Drive) [/dev/sdb] — Disk Utility'. The left sidebar shows 'Storage Devices' and 'Local Storage'. The main area displays the drive's serial number (07A1070890FFF8DC), world wide name (-), device (/dev/sdb), rotation rate (-), connection (USB at 480.0 Mb/s), and SMART status (Not Supported). A 'Safe Removal' button is present with the instruction 'Power down the drive so it can be removed'. Below this, a partition named 'linuxlive' is shown with a capacity of 4.0 GB FAT. The partition details are listed below:

Usage:	Filesystem
Partition Type:	W95 FAT32 (LBA) (0x0c)
Partition Flags:	Bootable
Type:	FAT (32-bit version)
Label:	linuxlive

Below the partition details, there are three buttons: 'Unmount Volume' (Unmount the volume), 'Check Filesystem' (Check and repair the filesystem), and 'Delete Partition' (Delete the partition). To the right of these buttons, there are two more buttons: 'Format Volume' (Erase or format the volume) and 'Edit Partition' (Change partition type, label and flags).

File system formats

By default, many USB flash disks are formatted in **FAT32**.

Other types are NTFS, ext4, ZFS.

FAT32 – max 4GB files

NTFS – maximum portability (also for use under windows)

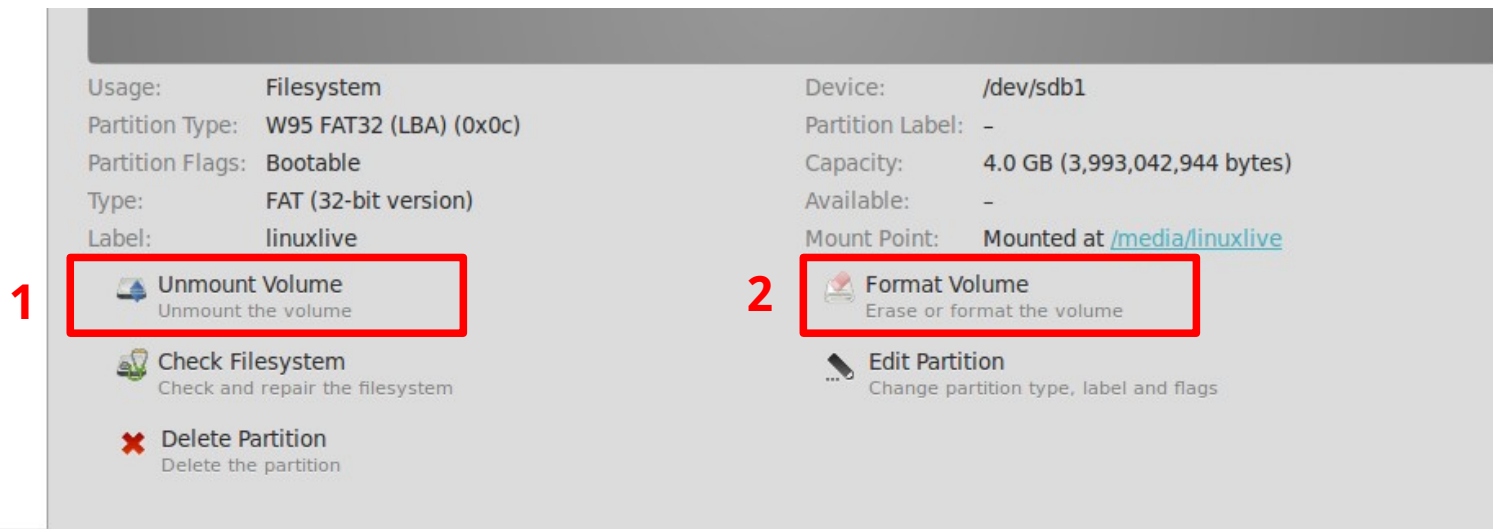
Ext4 – default file system in Linux,

Btrfs – the next default file system in Linux in the near future.

Example: formatting a USB disk

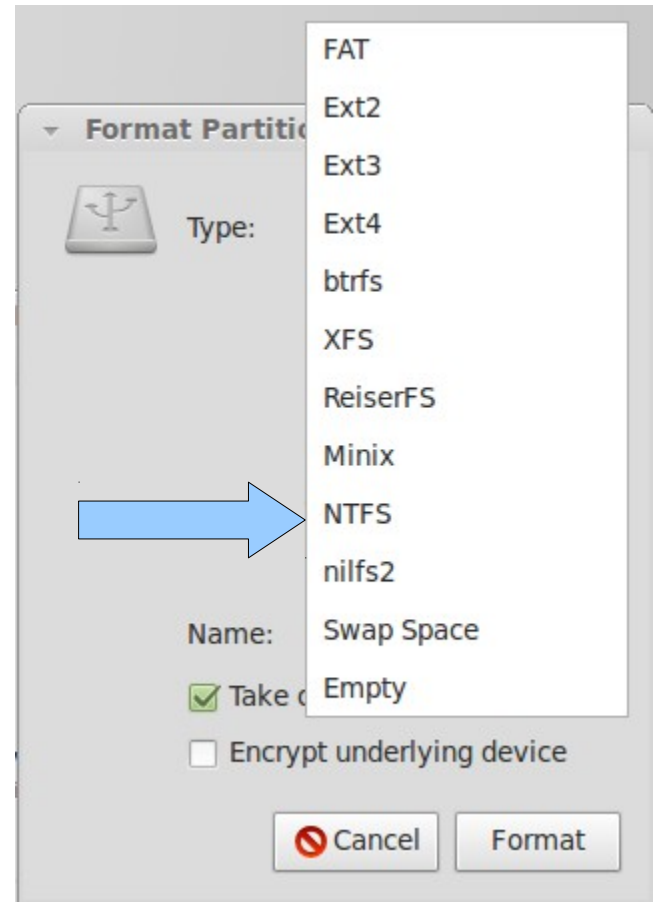
First **unmount** the device (red box 1 below). Now the USB disk can not be accessed anymore.

Next, choose **format** the device (red box 2 below).



Format disks with disk utility

Choose the type of file system you want to be on that device.



Format disks with disk utility

TDKMedia Trans-It Drive (TDKMedia Trans-It Drive) [/dev/sdb] — Disk Utility

File Help

Storage Devices

- Local Storage
joachim@localhost
 - PATA Host Adapter
82371AB/EB/MB PIIX4 IDE
 - CD/DVD Drive
VBOX VBOX CD-ROM
 - SATA Host Adapter
82801HM/HEM (IC...roller [AHCI mode])
 - 13 GB Hard Disk
ATA VBOX HARDDISK
- Peripheral Devices
USB, FireWire and other peripherals
- TDKMedia Trans-It Drive**
TDKMedia Trans-It Drive

Drive

Model: TDKMedia Trans-It Drive
Firmware Version: PMAP
Location: -
Write Cache: -
Capacity: 4.0 GB (3,997,171,712 bytes)
Partitioning: Master Boot Record

Serial Number: 07A1070890FFF8DC
World Wide Name: -
Device: /dev/sdb
Rotation Rate: -
Connection: USB at 480.0 Mb/s
SMART Status: ● Not Supported

Format Drive
Erase or partition the drive

Safe Removal
Power down the drive so it can be removed

Benchmark
Measure drive performance

Volumes

test
4.0 GB NTFS

Usage: Filesystem
Partition Type: W95 FAT32 (LBA) (0x0c)
Partition Flags: Bootable
Type: NTFS
Label: test

Device: /dev/sdb1
Partition Label: -
Capacity: 4.0 GB (3,993,042,944 bytes)
Available: -
Mount Point: Mounted at [/media/test](#)

Unmount Volume
Unmount the volume

Format Volume
Erase or format the volume

Check Filesystem
Check and repair the filesystem

Edit Partition
Change partition type, label and flags

Delete Partition
Delete the partition

Format disks with disk utility

The program disk-utility put a lot of commands at work behind the scenes.

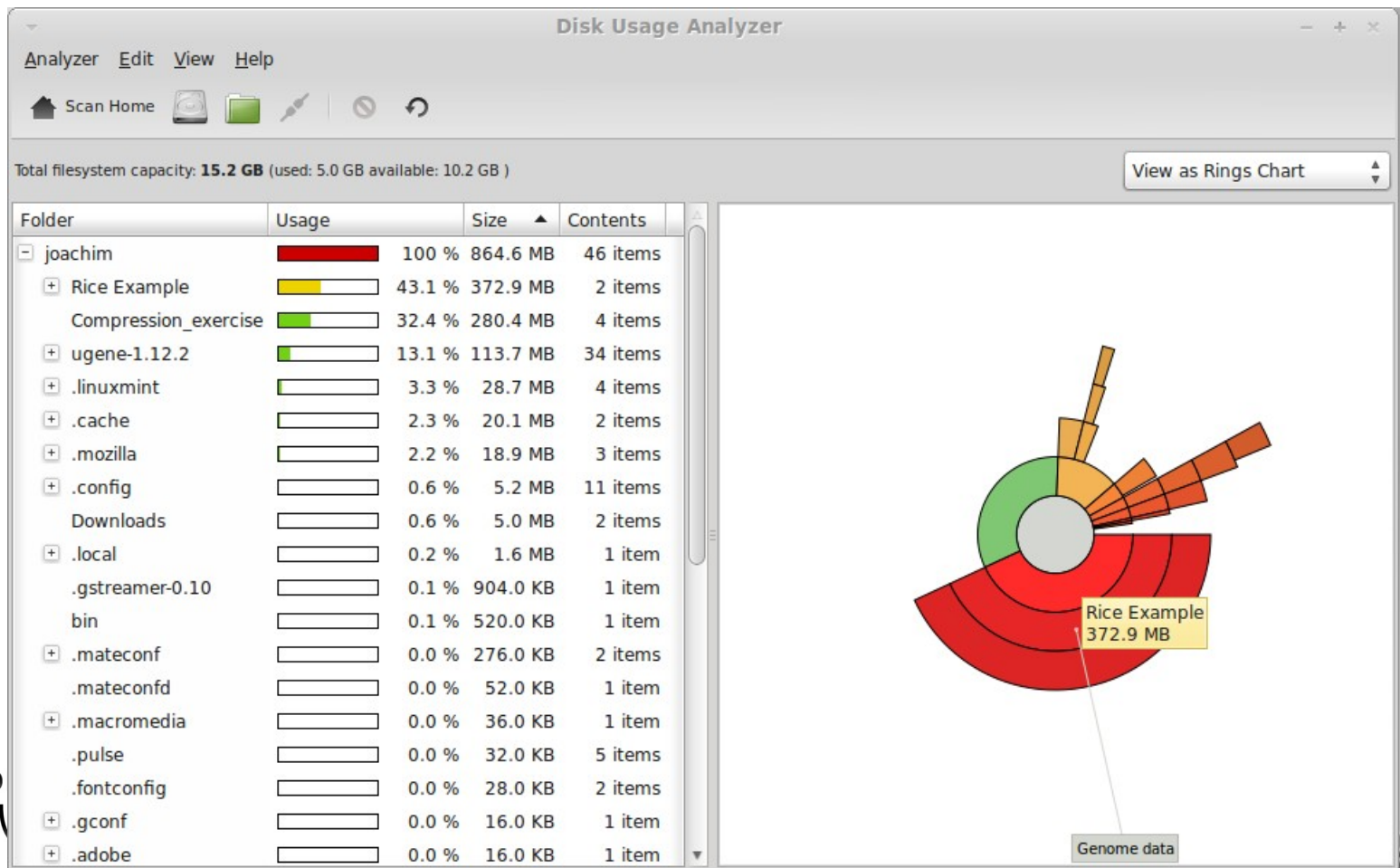
Some of the command line tools used:

- mount
- umount
- fdisk
- mkfs

You can read the man pages and search for guides on the internet if you want to get to know these (out of scope for this course).

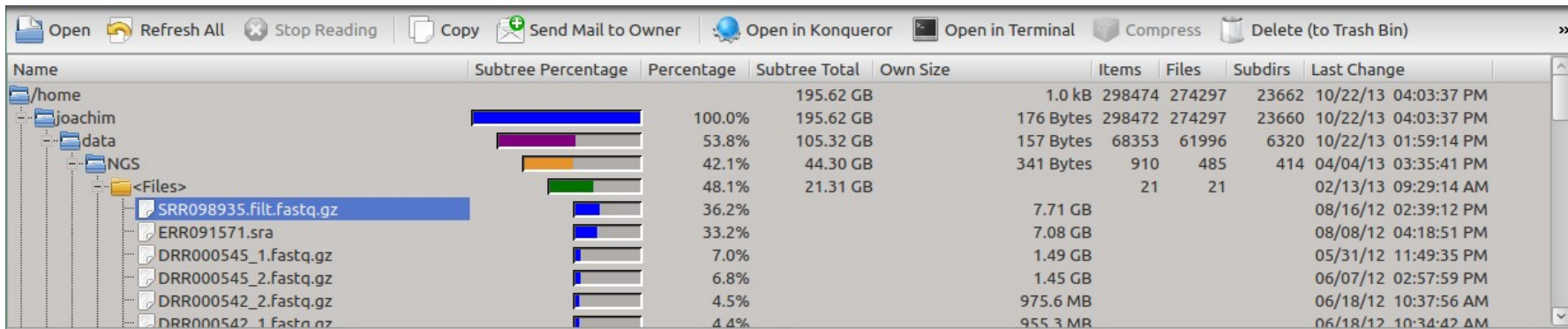
Checking storage space

By default 'disk usage analyzer'.



Checking storage space

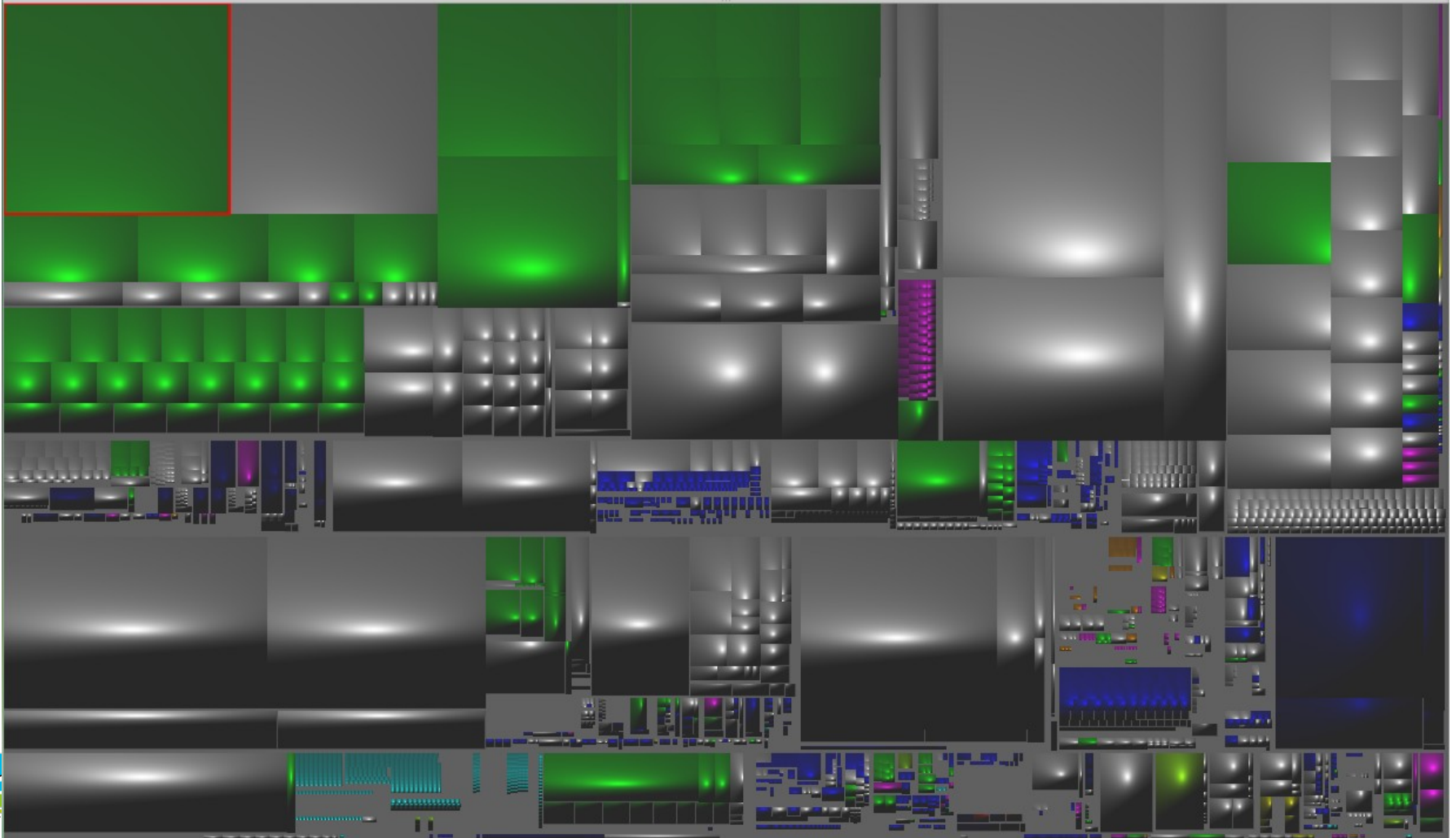
Bonus: K4DirStat. Not installed by default.



The image shows a screenshot of the K4DirStat application window. The window has a menu bar with options: Open, Refresh All, Stop Reading, Copy, Send Mail to Owner, Open in Konqueror, Open in Terminal, Compress, and Delete (to Trash Bin). Below the menu bar is a table with columns: Name, Subtree Percentage, Percentage, Subtree Total, Own Size, Items, Files, Subdirs, and Last Change. The table displays data for the /home directory and its subdirectories. The /home directory is highlighted in blue. The /home/data directory is highlighted in orange. The /home/data/NGS directory is highlighted in green. The /home/data/NGS directory contains several files, including SRR098935.filt.fastq.gz, ERR091571.sra, DRR000545_1.fastq.gz, DRR000545_2.fastq.gz, DRR000542_1.fastq.gz, and DRR000542_2.fastq.gz. The table shows that the /home directory is 195.62 GB in size, and the /home/data/NGS directory is 21.31 GB in size. The /home/data/NGS directory contains 21 files and 21 subdirectories. The /home/data/NGS directory is 48.1% of the total size of the /home directory.

Name	Subtree Percentage	Percentage	Subtree Total	Own Size	Items	Files	Subdirs	Last Change
/home			195.62 GB	1.0 kB	298474	274297	23662	10/22/13 04:03:37 PM
joachim			195.62 GB	176 Bytes	298472	274297	23660	10/22/13 04:03:37 PM
data			105.32 GB	157 Bytes	68353	61996	6320	10/22/13 01:59:14 PM
NGS			44.30 GB	341 Bytes	910	485	414	04/04/13 03:35:41 PM
<Files>			21.31 GB		21	21		02/13/13 09:29:14 AM
SRR098935.filt.fastq.gz		36.2%	7.71 GB					08/16/12 02:39:12 PM
ERR091571.sra		33.2%	7.08 GB					08/08/12 04:18:51 PM
DRR000545_1.fastq.gz		7.0%	1.49 GB					05/31/12 11:49:35 PM
DRR000545_2.fastq.gz		6.8%	1.45 GB					06/07/12 02:57:59 PM
DRR000542_1.fastq.gz		4.5%	975.6 MB					06/18/12 10:37:56 AM
DRR000542_2.fastq.gz		4.4%	955.3 MB					06/18/12 10:34:42 AM

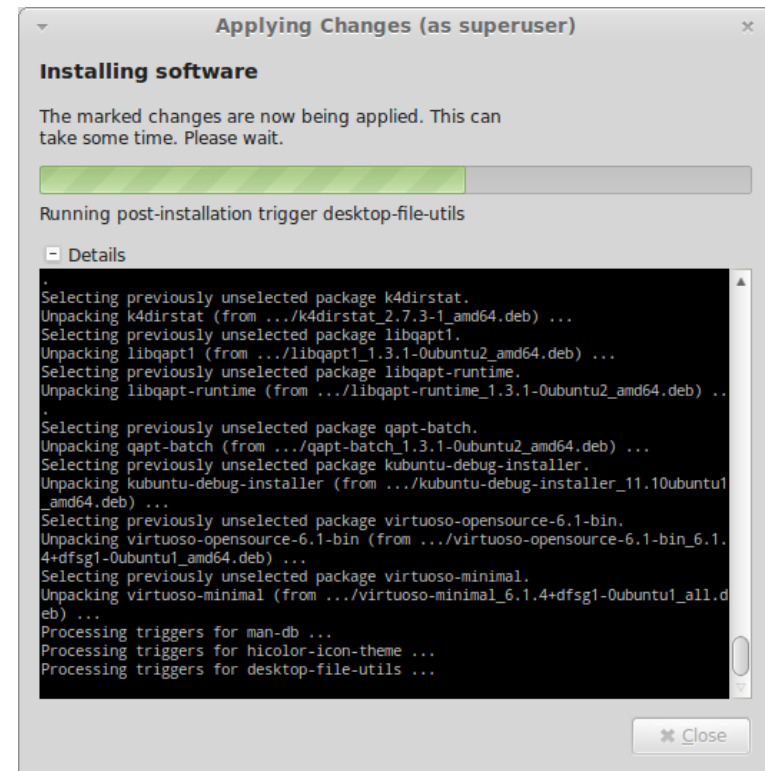
Name	Subtree Percentage	Percentage	Subtree Total	Own Size	Items	Files	Subdirs	Last Change
/home			195.62 GB		1.0 kB	298474	274297	23662 10/22/13 04:03:37 PM
joachim								
data								
NGS								
<Files>								
SRR098935.filt.fastq.gz								
ERR091571.sra								
DRR000545_1.fastq.gz								
DRR000545_2.fastq.gz								
DRR000542_2.fastq.gz								
DRR000542_1.fastq.gz								



K4Dirstat is a KDE package

Rehearsal: what is KDE?

Bonus: what happens when you install this package on our system?



Space left on disks with df

To check the storage that is used on the different disks: **df -h**

```
~/ $ df -h
Filesystem      Size  Used Avail Use% Mounted on
/dev/sda1       12G   5.3G   5.7G   49% /
udev           490M   4.0K   490M    1% /dev
tmpfs           200M   920K   199M    1% /run
none            5.0M     0    5.0M    0% /run/lock
none            498M   76K   498M    1% /run/shm
/dev/sdb1       3.8G   20M   3.7G    1% /media/test

~/ $ df -h .
```

The size of directories

To check the size of files or directories: **du**

```
~/ $ du -sh *  
520K  bin  
281M  Compression_exercise  
4.0K  Desktop  
4.0K  Documents  
5.0M  Downloads  
4.0K  Music  
4.0K  Pictures  
4.0K  Public  
373M  Rice Example  
4.0K  Templates  
4.0K  test  
17Mtest.img  
114M  ugene-1.12.2  
4.0K  Videos
```

* means 'anything'.

This is called 'globbing':
* is a wild card symbol.

Wildcards on the command line

Wildcards are used to describe the names of **files/dirs**.

*

On that position, any length of string is allowed
e.g. **s*** matches: **san**, **sdd**, **sanitisation**, **sam.alignment**,...

?

On that position, any character is allowed.
e.g. **saniti?ation** matches: **sanitisation**, **sanitiration**, ...

[]

On that position, the character may be one of the characters between [],

e.g. **saniti[sz]ation** matches: **sanitisation** and **sanitization**

Wildcards on the command line

Many tools that require an **argument** to point to files or directories accept these wildcards.

```
~/ $ du -sh Do*
```

Wildcards on the command line

Many tools that require an **argument** to point to files or directories accept these wildcards.

```
~/ $ du -sh Do*  
4.0K Documents  
20G Downloads
```

Wildcards on the command line

Many tools that require an **argument** to point to files or directories accept these wildcards.

```
~/ $ ls *.fastq
```

Wildcards on the command line

Many tools that require an **argument** to point to files or directories accept these wildcards.

```
~/ $ ls *.fastq
ERR148552_1.fastq      ERR148552_2.fastq
testout.fastq
ERR148552_1_prinseq_good_zzwI.fastq  test.fastq
```

Keywords

Compression

Archive

Symbolic link

mounting

File system format

partition

Recursively

df

du

unlink

Break

