

BASIC UNIX

Module

Introduction

Unix is the standard operating system on most large computer systems in scientific research, in the same way that Microsoft Windows is the dominant operating system on desktop PCs.

Unix and MS Windows both perform the important job of managing the computer's hardware (screen, keyboard, mouse, hard disks, network connections, etc) on your behalf. They also provide you with tools to manage your files and to run application software.

They both offer a graphical user interface. On Unix systems, this is called the X Window System, or just X.

Unix is a powerful, robust and stable operating system which allows dozens of people to run programs on the same computer at the same time. This is why it is the preferred operating system for large-scale scientific computing. It runs on all kinds of machines, from desktop PCs to supercomputers.

Aims

The aim of this module is to introduce UNIX and cover some of the basics that will allow you to run some of the programs used in this workshop. Several of the programs that you are going to use during the workshop, plus many others that are useful for bioinformatics analyses, are run in UNIX. This module is only designed to provide a very brief introduction to some of the features and useful commands of UNIX.

During this module we will also obtain a genome sequence that will be used in the next module, and examine the basic structure of an EMBL entry.

Why use UNIX?

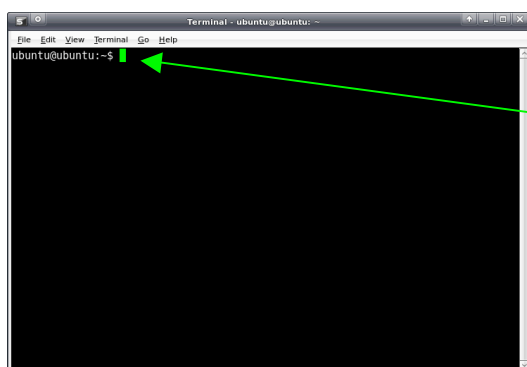
- UNIX is a well established, very widespread operating system.
- Command line driven, with a huge number of often terse, but powerful commands
- In contrast to Windows, it is designed to allow many users to run their programs simultaneously on the same computer
- Designed to work in computer networks - for example, most of the Internet is UNIX based
- It is used on many of the powerful computers at bioinformatics centres and also on many workstations
- UNIX is not a monolithic entity. There are numerous different UNIX operating systems. Some of them are freely distributed such as Linux which was originally created to provide a free UNIX on personal PCs. This operating system is now so popular that it has been ported to many different system architectures.

Getting started

In this workshop, we will be using desktop PCs which run Linux, a version of UNIX which was specially designed for PCs.

We will use a terminal window to type in our UNIX command line. This is similar to the "Command Prompt" window on MS Windows systems, which allows the user to type DOS commands to manage files.

You should see a window labelled "Terminal" which will be empty except for a '\$' character at the top left. The '\$' character is the UNIX prompt, similar to "C:\\" in DOS. Note: the prompt will often be different on different UNIX computers, for example it may be displayed as a '%' character.

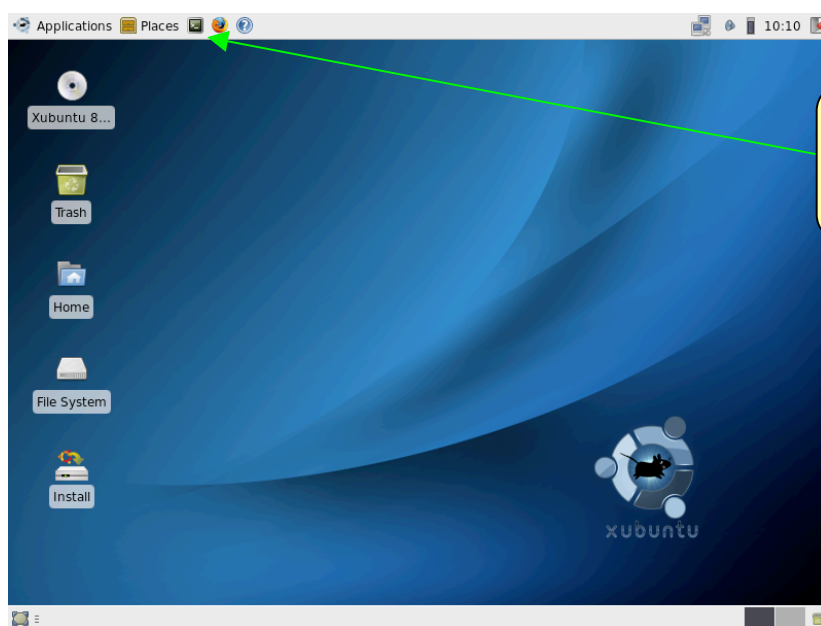


UNIX prompt

You can type commands directly into the terminal at the '\$' prompt.

A list of useful commands can be found on the next page.

Many of them are two- or three-letter abbreviations. The earliest UNIX systems (*circa* 1970) only had slow Teletype terminals, so it was faster to type 'rm' to remove a file than 'delete' or 'erase'. This terseness is a feature of UNIX which still survives.



To get a terminal window click on the terminal icon

The command line

All UNIX programs may be run by typing commands at the UNIX prompt **\$**. The command line tells the computer what to do.

You may subtly alter these commands by specifying certain options when typing in the command line.

Command line Arguments

Typing any of the commands listed above at the UNIX prompt with the appropriate variables such as files names or directories will result in the tasks being performed on pressing the enter key.

Additional arguments or flags can be added to the commands to affect the way the command works. For example:

The **cal** command prints a calendar for a month or a year

If you type in just **cal** with no month or year, you get the calendar for the current month

If you type **cal** and a year you get the calendar for that year

```
$ cal 2000 [enter]
```

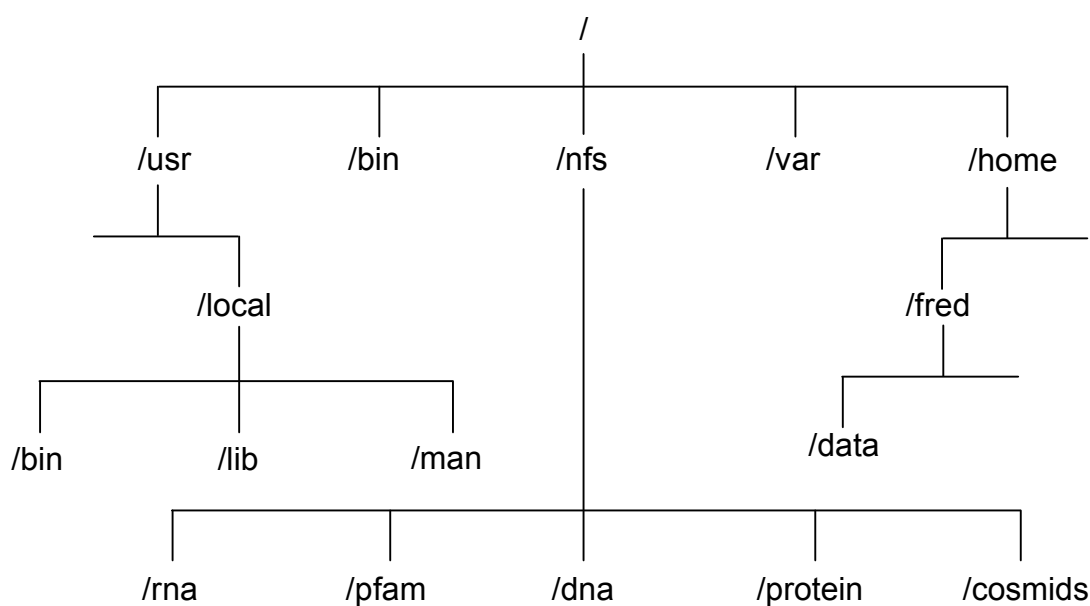
Additional arguments for the commands are not covered here, but if you want to find out what arguments are available, or want to find out more about a UNIX command, type **man** followed by the UNIX command

```
$ man cal [enter]
```

Files and Directories

Directories are the UNIX equivalent of folders on a PC or Mac. They are organised in a hierarchy, so directories can have sub-directories and so on. Directories are very useful for organising your work and keeping your account tidy - for example, if you have more than one project, you can organise the files for each one into different directories to keep them separate. You can think of directories as rooms in a house. You can only be in one room (directory) at a time. When you are in a room you can see everything in that room easily. To see things in other rooms, you have to go to the appropriate door and crane your head around. UNIX works in a similar manner, moving from directory to directory to access files. The location or directory that you are in is referred to as the current working directory.

Directory structure example



Therefore if there is a file called `genome.seq` in the **dna** directory its location or full pathname can be expressed as `/nfs/dna/genome.seq`.

For the actual directory structure you will be using during the workshop, see Appendix IV.

General Points

UNIX is pretty straightforward, but there are some general points to remember that will make your life easier:

UNIX is case sensitive - typing "ls" is not the same as typing "LS".

You need to put a space between a command and its argument - for example, "more myfile" will show you the contents of the file called myfile; "moremyfile" will just give you an error!

UNIX is not psychic! If you mis-spell the name of a command or the name of a file, it will not understand you.

Many of the commands are only a few letters long; this can be confusing until you start to think logically about why those letters were chosen - ls for list, rm for remove and so on.

Often when you have problems with UNIX, it is due to a spelling mistake, or perhaps you have omitted a space.

If you want to know more about UNIX and its commands there are plenty of resources available that provide a more comprehensive guide:-

(e.g. <http://unixhelp.ed.ac.uk> or <http://unix.t-a-y-l-o-r.com/>).

In what follows, we shall use the following typographical conventions:

Characters written in **bold typewriter font** are commands to be typed into the computer as they stand.

Characters written in *italic typewriter font* indicate non-specific file or directory names.

Words inserted within square brackets [Ctrl] indicate keys to be pressed.

So, for example,

\$ ls anydirectory [Enter]

means "at the UNIX prompt \$, type ls followed by the name of some directory, then press the key marked Enter"

Don't forget to press the [Enter] key: commands are not sent to the computer until this is done.

Some useful UNIX commands

| Command | What it does |
|--------------|--|
| ls | Lists the contents of the current directory |
| mkdir | Makes a new directory |
| mv | Moves or renames a file |
| cp | Copies a file |
| rm | Removes a file |
| cat | Concatenates files |
| more | Displays the contents of a file one page at a time |
| head | Displays the first ten lines of a file |
| tail | Displays the last ten lines of a file |
| cd | Changes current working directory |
| pwd | Prints working directory |
| find | Finds files matching an expression |
| grep | Searches a file for patterns |
| wc | Counts the lines, words, characters, and bytes in a file |
| kill | Stops a process |
| jobs | Lists the processes that are running |

Exercise

The following exercise introduces a few useful UNIX commands and provides examples of how they can be used.

Many people panic when they are confronted with an UNIX prompt! Don't! The exercise is designed to be step-by-step, so all the commands you need are provided in the text. If you get lost ask a demonstrator. If you are a person skilled at UNIX, be patient it is only a short exercise.

Finding where you are and what you've got

pwd Print the working directory

As seen previously directories are arranged in a hierarchical structure. To determine where you are in the hierarchy you can use the **pwd** command to display the name of the current working directory. The current working directory may be thought of as the directory you are in, i.e. your current position in the file-system tree

To find out where you are type

```
pwd [enter]
```

You will see that you are in your home directory. We need to move into the UNIX directory

UNIX is case sensitive, **PWD** is not the same as **pwd**

cd Change current working directory

The **cd** command will change the current working directory to another, in other words allow you to move up or down in the directory hierarchy. First of all we are going to move into the UNIX directory below. To do this type:

```
cd UNIX [enter]
```

Now use the **pwd** command to check your location in the directory hierarchy.

ls List the contents of a directory

To find out what are the contents of the current directory type

```
ls [enter]
```

The **ls** command lists the contents of your current directory, this includes files and directories. You should see that there are 4 files called:
AL513382.embl, MAL13P1.dna, MAL13P1.tab, Malaria.fasta

Changing and moving what you've got

cp Copy a file

`cp file1 file2` is the the command which makes a copy of file1 in the current working directory and calls it file2

What you are going to do is make a copy of `AL513382.embl`. This file contains the genome of *Salmonella typhi* strain CT18 in EMBL format. The new file will be called `S_typhi.embl`

```
cp AL513382.embl S_typhi.embl [enter]
```

If you use the `ls` command to check the contents of the current directory you will see that there are now two files, `AL513382.embl` and `S_typhi.embl`.

rm Delete a file

This command removes a file permanently so take care!

You are now going to remove the old version of *S. typhi* genome file, `AL513382.embl`

```
rm AL513382.embl [enter]
```

The file will be removed. Use the `ls` command to check the contents of the current directory to see that `AL513382.embl` has been removed.

UNIX as a general rule does exactly what you ask, and does not ask for confirmation. Unfortunately there is no "recycle bin" on the command line to recover the file from, so you have to be careful.

cd Change current working directory

As before the `cd` command will change the current working directory to another, in other words allow you to move up or down in the directory hierarchy. First of all we are going to move into the directory above, type:

```
cd .. [enter]
```

Now use the `pwd` command to check your location in the directory hierarchy.

Next, we are going to move into the `Module_1_Artemis` directory.

To change to the `Module_1_Artemis` directory type:

```
cd Module_1_Artemis [enter]
```

use the **ls** command to check the contents of the directory

mv Move a file

To move a file from one place to another use the **mv** command. This moves the file rather than copies it, therefore you end up with only one file rather than two.

When using the command the path or pathname is used to tell UNIX where to find the file. You refer to files in other directories by using the list of hierarchical names separated by slashes. For example, the file *bases* in the directory *genome* has the path **genome/bases**

If no path is specified UNIX assumes that the file is in the current working directory.

What you are going to do is move the file `S_typhi.embl` from the UNIX directory, to the current working directory

```
mv ../UNIX/S_typhi.embl .        [enter]
```

Use the **ls** command to check the contents of the current directory to see that `S_typhi.embl` has been moved.

`../UNIX/S_typhi.embl` specifies that `S_typhi.embl` is in the UNIX directory. If the file was in the directory above, the path would change to: `../S_typhi.embl`

- `.` specifies the current working directory
- `..` specifies the directory above the current working directory

The command can also be used to rename a file in the current working directory. Previously we used the **cp** command, but **mv** provides an alternative without the need to delete the original file. Therefore we could have used:

```
mv AL513382.embl S_typhi.embl [enter]
```

Viewing what you've got

more Display file contents

This command displays the contents of a specified file one screen at a time.

You are now going to look at the contents of `S_typhi.embl`.

```
more S_typhi.embl [enter]
```

The contents of `S_typhi.embl` will be displayed one screen at a time, to view the next screen press the **space bar**. The percentage of the file that has been viewed so far will be displayed at the bottom of the screen. As `S_typhi.embl` is a large file this will take a while, therefore you may want to escape or exit from his command. To do this press the **control** and **c** keys simultaneously, this kills the **more** command, and returns you to the UNIX prompt. **more** can also scroll backwards if you hit the **b** key. Another useful feature is the **slash key**, `/`, to search for an expression in the file.

head Display the first ten lines of a file

tail Display the last ten lines of a file

Sometimes you may just want to view the text at the beginning or the end of a file, without having to display all of the file. The `head` and `tail` commands can be used to do this.

You are now going to look at the beginning of `S_typhi.embl`.

```
head S_typhi.embl [enter]
```

To look at the end of `S_typhi.embl` type:

```
tail S_typhi.embl [enter]
```

The amount of the file that is displayed can be increased by adding extra arguments. To increase the number of lines viewed from 10 to 100 add the `-100` argument to the command. For example to view the last 100 lines of `S_typhi.embl` type:

```
tail -100 S_typhi.embl [enter]
```

Do this for both `head` and `tail` commands. What type of information is at the beginning and end of the EMBL format file?

cat Join files together

Having looked at the beginning and end of the `S_typhi.embl` file you should notice that in EMBL format files the annotation comes first, then the DNA sequence at the end.

If you had two separate files containing the annotation and the DNA sequence, both in EMBL format, it is possible to concatenate or join the two together to make a single file like the `S_typhi.embl` file you have just looked at. The UNIX command **cat** can be used to join two or more files into a single file. The order in which the files are joined is determined by the order in which they appear in the command line.

For example, we have two separate files, `MAL13P1.dna` and `MAL13P1.tab`, that contain the DNA and annotation, respectively, from the *P. falciparum* genome.

By typing the command line:

```
cat MAL13P1.tab MAL13P1.dna > MAL13P1.embl [enter]
```

`MAL13P1.tab` and `MAL13P1.dna` will be joined together and written to a file called `MAL13P1.embl`

The `>` symbol in the command line directs the output of the **cat** program to the designated file `MAL13P1.embl`

wc Counts the lines, words or characters

By typing the command line:

```
ls | wc -l [enter]
```

The above command uses **wc** to count the number of files that are listed by **ls**.

The | symbol in the command line connects the two commands into a single operation for simplicity. You can connect as many commands as you want:

```
$ ls | grep ".embl" | wc -l
```

grep Searches a file for patterns

grep is a powerful tool to search for patterns in a file.

In the examples below, we are going to use the file called *Malaria.fasta* that contains the set of *P. falciparum* chromosomes in FASTA format. A FASTA file has the following format:

```
>Sequence Header  
CTAAACCTAAACCTAAACCCTGAACCCTAA...  
...
```

Therefore if we want to get the sequence headers, we can extract the lines that match the > symbol:

```
grep '>' Malaria.fasta [enter]
```

By typing the command line:

```
grep -c '>' Malaria.fasta [enter]
```

The > symbol is placed in quotes as this stops the shell from interpreting the > as an instruction for where to put the output.

The -c option prints only a count of matching lines. Therefore in this example we will display the number of sequence entries that this file contains.

find Finds files matching an expression

The **find** command is similar to **ls** but in many ways it is more powerful. It can be used to recursively search the directory tree for a specified path name, seeking files that match a given Boolean expression (a test which returns true or false)

find . -name "*.embl"

This command will return the files which name has the .embl suffix.

find . -type d

This command will return all the subdirectories contained in the current directory.

These is just two basic examples but it is possible to search in many other ways:

| | |
|--------|-------------------------------------|
| -mtime | search files by modifying date |
| -atime | search files by last access date |
| -size | search files by file size |
| -user | search files by user they belong to |

You need to be careful with quoting when using wildcards.

The wildcard * symbol represents a string of any character and of any length.

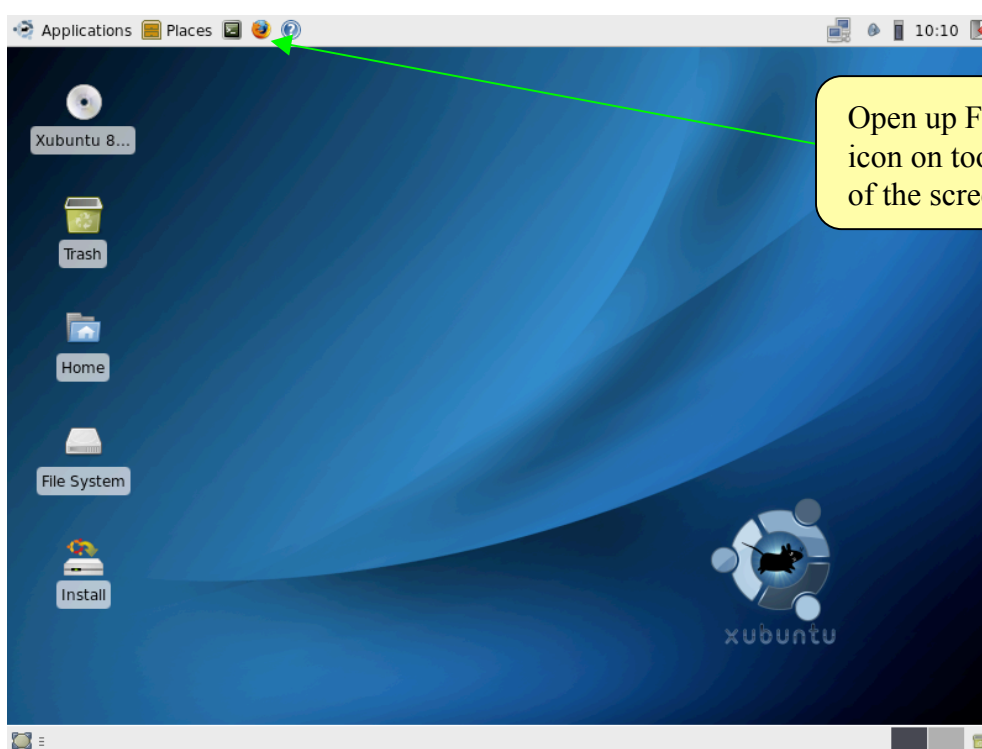
Obtaining and transferring information

The first step in exploiting genome sequences is obtaining your genome sequence. As time goes by there are more and more genome sequences available, from an ever increasing number of locations. Typically a complete genome sequence project is quite large, and therefore the files containing the data are going to be quite unwieldy. One of the simplest ways in which such information can be obtained is using **ftp** or 'file transfer protocol'. This a method of transferring information from a remote machine to the computer you are working on.

The **ftp** command can be used in UNIX to connect to a remote machine specified in the command line. Once a connection is established it is possible to both send (upload) and receive (download) data. However as we are limited for time we will not use this method, and instead use a more user-friendly method.

Using ftp on the internet

In addition to UNIX, ftp is also available on most Macs and PCs and allows you to transfer files readily between different computers worldwide. It is worth learning how to use **ftp**; most machines will have a graphical **ftp** interface and this makes file transfer very easy. Unfortunately there are a large number of alternatives and we can't show them all to you. Instead, we'll use **ftp** to download information to the current working directory on the computer you are working using the Firefox web browser. You are now going to an **ftp** web page where you are going to download the DNA sequence for the *Salmonella typhi* genome.



In the location box type in the address http://www.sanger.ac.uk/Projects/S_typhi/ and press return. The page below should appear.

Salmonella typhi

Sequence and annotation of *Salmonella enterica* serovar Typhi (*S. typhi*) strain CT18

Image source: [GlaxoSmithKline Biologicals, Belgium](#)

The annotation for the *S. typhi* CT18 genome is now accessible through [GeneDB](#). GeneDB allows searching, browsing and download of the full annotation.

The sequenced strain, CT18, is a highly pathogenic, multiple drug resistant strain isolated from a typhoid patient in Cho Quan Hospital, Ho Chi Minh City, Vietnam. The chromosome sequence is 4,809,037 bp in length with a G+C content of 52.09%, and was generated from 97,000 shotgun reads. The start of the sequence was chosen to correspond with the start of the published [E. coli K12 sequence](#). There are 4,599 protein-coding genes (including 204 pseudogenes). Both the sequence and annotation have been deposited in the public databases with the accession number [AL513382](#). In addition, there are two plasmids: pHCM1 (218,150 bp, 249 CDS) with the accession number [AL513383](#), and pHCM2 (106,516 bp, 131 CDS) with the accession number [AL513384](#).

The published sequence data of strain CT18, for the chromosomal sequence and both plasmids, can be searched using our [BLAST server](#) and are also available for download from our [ftp site](#), in fasta format, [Artemis](#) format and as translated peptide sequence.

The sequence and annotation was produced in collaboration with [Gordon Dougan](#) of the [Department of Biochemistry, Imperial College, London](#).

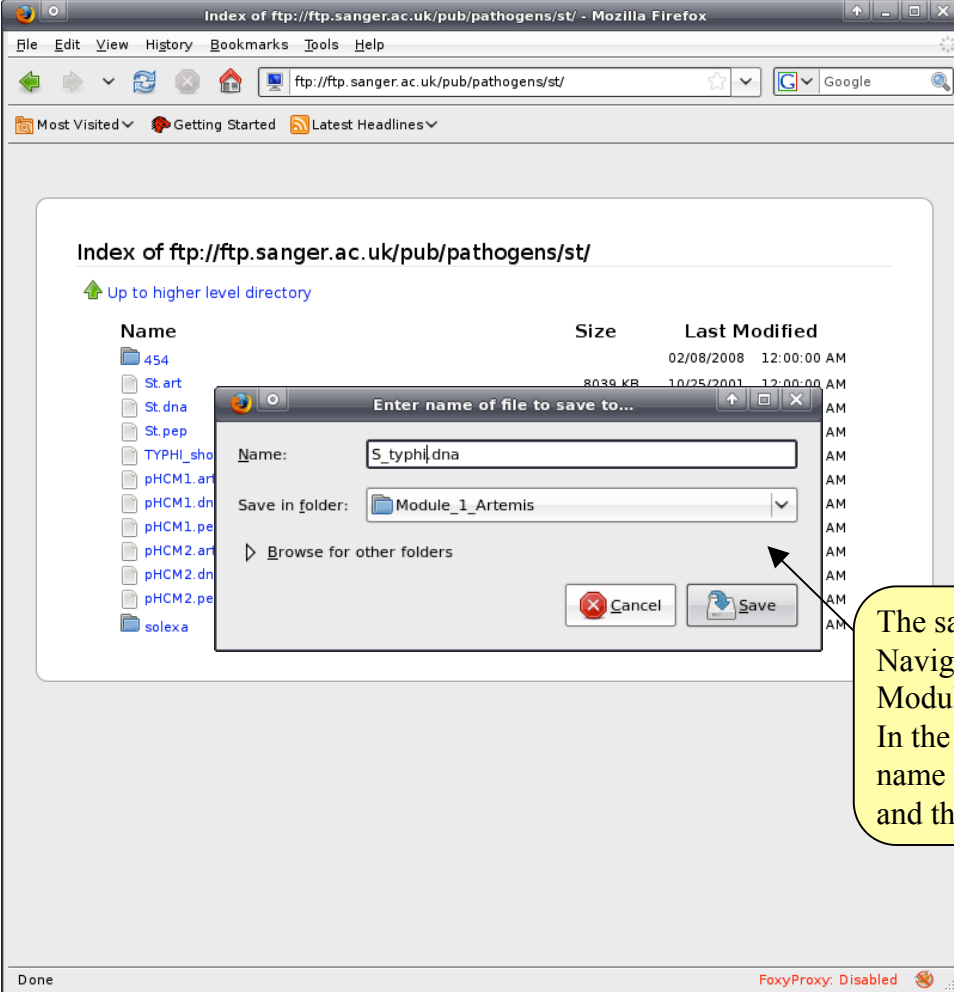
Click on the ftp access link

Index of <ftp://ftp.sanger.ac.uk/pub/pathogens/st/>

Up to higher level directory

| Name | Size | Last Modified |
|-------------------|----------|------------------------|
| 454 | | 02/08/2008 12:00:00 AM |
| St.art | 8020 KB | 10/25/2001 12:00:00 AM |
| St.dna | 4775 KB | 10/25/2001 12:00:00 AM |
| St.pep | 1658 KB | 10/25/2001 12:00:00 AM |
| TYPHI_shotgun.dbs | 47684 KB | 04/03/2008 12:00:00 AM |
| pHCM1.art | 426 KB | 10/25/2001 12:00:00 AM |
| pHCM1.dna | 217 KB | 10/25/2001 12:00:00 AM |
| pHCM1.pep | 74 KB | 10/25/2001 12:00:00 AM |
| pHCM2.art | 240 KB | 10/25/2001 12:00:00 AM |
| pHCM2.dna | 106 KB | 10/25/2001 12:00:00 AM |
| pHCM2.pep | 39 KB | 10/25/2001 12:00:00 AM |
| solexa | | 06/30/2008 09:40:00 AM |

Right click on the St.dna file and choose 'Save Link As'.



Index of ftp://ftp.sanger.ac.uk/pub/pathogens/st/

Up to higher level directory

| Name | Size | Last Modified |
|-----------|---------|------------------------|
| 454 | | 02/08/2008 12:00:00 AM |
| St.art | | |
| St.dna | 8039 KB | 10/25/2001 12:00:00 AM |
| St.pep | | |
| TYPHI_sho | | |
| pHCM1.ar | | |
| pHCM1.dn | | |
| pHCM1.pe | | |
| pHCM2.ar | | |
| pHCM2.dn | | |
| pHCM2.pe | | |
| solexa | | |

Enter name of file to save to...

Name: S_typhi.dna

Save in folder: Module_1_Artemis

Browse for other folders

Cancel Save

The save menu will appear. Navigate to the Module_1_Artemis directory. In the selection box type the name of the file as S_typhi.dna and then press OK

The file containing the DNA sequence for the genome of *S. typhi* will now be saved to the Module_1_Artemis directory

Database Entries

The `S_typhi.embl` file we previously looked at was obtained from the EMBL database at the European Bioinformatics Institute (<http://srs.ebi.ac.uk/>) and is presented in a specific format with a series of defined qualifiers and keys (see below and Appendix III) to help identify the different components of an entry.

Below is an example of a small EMBL entry with the different features of the entry highlighted

```

ID AF060869 standard; DNA; PRO; 27290 BP.
XX
AC AF060869;
XX
SV AF060869.1
XX
DT 17-JUL-1998 (Rel. 56, Created)
DT 17-JUL-1998 (Rel. 56, Last updated, Version 1)
XX
DE Salmonella typhimurium excision nuclease UvrA (uvrA) gene, partial cds;
DE single-strand binding protein (ssb) gene, complete cds; tRNA-Thr gene,
DE complete sequence; pathogenicity island SPI-4 operon, complete sequence;
DE yjcB gene, complete cds; and yjcC gene, partial cds.
OS Salmonella typhimurium
OC Bacteria; Proteobacteria; gamma subdivision; Enterobacteriaceae;
OC Salmonella.
XX
RN [1]
RP 1-27290
RX MEDLINE; 98298059.
RA Wong K.K., McClelland M., Stillwell L.C., Sisk E.C., Thurston S.J.,
RA Saffer J.D.;
RT "Identification and sequence analysis of a 27-kilobase chromosomal fragment
RT containing a Salmonella pathogenicity island located at 92 minutes on the
RT chromosome map of Salmonella enterica serovar typhimurium LT2";
RL Infect. Immun. 66(7):3365-3371(1998).
DR SPTREMBL; 085309; 085309.
DR SPTREMBL; 085310; 085310.
XX
FH Key Location/Qualifiers
FH
FT source 1..27290
FT /db_xref="taxon:602"
FT /organism="Salmonella typhimurium"
FT /strain="LT2"
FT /map="92 minutes"
FT CDS complement(1..312)
FT /codon_start=1
FT /db_xref="SPTREMBL:085309"
FT /transl_table=11
FT /gene="uvrA"
FT /product="excision nuclease UvrA"
FT /protein_id="AAC26637.1"
FT /translation="MDKIEVRGARTHNLKNINRVIPRDKLIVVTGLSGSGKSSLAFDTL
FT YAEQGRRYVESLSAYARQFLSLMEKPDVDHIEGLSPAISIEQKSTSHNPRSTVGTITEI
FT "
FT CDS 583..1083
FT /codon_start=1
FT /db_xref="SPTREMBL:085310"
FT /transl_table=11
FT /gene="ssb"
FT /product="single-strand binding protein"
FT /protein_id="AAC26638.1"
FT /translation="MILVGNPQDPVRYMPSGGAVANLTLATSESWPKQTGEMKEQT
FT EIHRRVVMFKLAIVAGEYLRLLSSQVYIEGLRLTRKRTDQNCQERYTTELTSDRRVMQI
FT LGGPKGGGAPAGGHNRGLGSPQQPQQPGGNQFNNGAQRPPQSAFAPSKPPMDFDD
FT IPF"
FT tRNA 1898..1970
FT /note="putative"
FT /product="tRNA-Thr"
XX
SQ Sequence 27290 BP: 7965 A: 5530 C: 6661 G: 7134 T: 0 other:
gatctcgata atagtaccoc cogtagagcg cggcttgtag gatgctgatt totgttcoaat 60
tgagatcgcg ggcgatagcc cctcaatcag gtcgacatcc ggtttttcca tgasgcacaa 120
aaactcgcgc gctgaagcgc agagcagatc aacgtaacga cgcctgacct cggcatacag 180
ttgtttatc gtatcatgac aggctgagaa gcttttcaga agaggacact tataaaataa 2520
aggcttggtt agaagacaaa atcaatagta atttattgat agaaatggtt attcctcagg 2580
//

```

Key

Qualifier

EMBL Header

Annotation

Sequence

Two-character line code indicates the type of information contained in the line

In addition to the EMBL database, there are the mirror databases, Genbank (NCBI) and DDBJ (National Institute of Genetics, Japan), which contain the same sequence entries, but have slight differences in the way in which the information is presented. The next two pages contain the text of the complete entries for same sequence from the EMBL and GenBank databases, compare the two entries and identify the differences.

EMBL Entry

```

ID  ECRSMA      standard; DNA; PRO; 500 BP.
XX
AC  L40173;
XX
SV  L40173.1
XX
DT  10-AUG-1995 (Rel. 44, Created)
DT  04-MAR-2000 (Rel. 63, Last updated, Version 4)
XX
DE  Erwinia carotovora repressor (rsmA) gene, complete cds.
XX
KW  repressor; rsmA gene.
XX
OS  Pectobacterium carotovorum
OC  Bacteria; Proteobacteria; Gammaproteobacteria; Enterobacteriaceae;
OC  Pectobacterium.
XX
RN  [1]
RP  1-500
RA  Cui Y., Chatterjee A., Liu Y., Dumenyo C.K., Chatterjee A.K.;
RT  "Identification of a global repressor gene, rsmA, of Erwinia carotovora
RT  subsp. carotovora that controls extracellular enzymes,
RT  N-(3-oxohexanoyl)-L-homoserine lactone, and pathogenicity in soft-rotting
RT  Erwinia spp";
RL  J. Bacteriol. 177(17):0-0(1995).
XX
DR  GOA; Q47620; Q47620.
DR  SWISS-PROT; Q47620; CSRA_ERWCA.
XX
FH  Key          Location/Qualifiers
FH
FT  source       1..500
FT              /db_xref="taxon:554"
FT              /organism="Pectobacterium carotovorum"
FT              /strain="71"
FT              /sub_species="carotovora"
FT              /gene="rsmA"
FT  CDS          246..431
FT              /codon_start=1
FT              /db_xref="GOA:Q47620"
FT              /db_xref="SWISS-PROT:Q47620"
FT              /note="putative"
FT              /transl_table=11
FT              /gene="rsmA"
FT              /function="global repressor"
FT              /protein_id="AAA74502.1"
FT              /translation="MLILTRRVGETLIIGDEVTVTVLVGKGNQVRIGVNPKEVSVHRE
FT              EIYQRIQAEKSQPTSY"
XX
SQ  Sequence 500 BP; 140 A; 101 C; 120 G; 139 T; 0 other;
    ggatccggca agcaggatag aaagtgtggtt accttcagat attctgaagc ttacatgct      60
    cagttctggtt gttgtgataa caaaagcaca agctactgat atcgactaaa ctaacaagta      120
    gtgacaaacc ggagtgtgat ggtgtgggta taccatcgtc taggtttacg ttttcacagc      180
    acatgatgga taatggcggg gagacagaga gacccgactc tttataatct ttcaaggagc      240
    aaagaatgct tattttgact cgtcagattg gcgaaaccct catcatcggc gatgaggtaa      300
    cggttaccgt attaggagtg aaaggcaacc aggtgcgtat tgggtgtaat gcacctaaag      360
    aggtttctgt ccaccgtgaa gagatctatc agcgtattca ggccgaaaaa tctcaaccaa      420
    cgtcatattg attgacaatg cgtctcgtgt tcgcggggacg caattgttat ttccgggttt      480
    tccccacac  atttatcgat                                     500
//

```

GenBank Entry

```

LOCUS      ERWRSMA                      500 bp    DNA     linear   BCT 19-AUG-1995
DEFINITION Erwinia carotovora repressor (rsmA) gene, complete cds.
ACCESSION  L40173
VERSION    L40173.1  GI:927031
KEYWORDS   repressor; rsmA gene.
SOURCE     Pectobacterium carotovorum
  ORGANISM Pectobacterium carotovorum
            Bacteria; Proteobacteria; Gammaproteobacteria; Enterobacteriaceae;
            Pectobacterium.
REFERENCE  1 (bases 1 to 500)
  AUTHORS  Cui,Y., Chatterjee,A., Liu,Y., Dumenyo,C.K. and Chatterjee,A.K.
  TITLE    Identification of a global repressor gene, rsmA, of Erwinia
            carotovora subsp. carotovora that controls extracellular enzymes,
            N-(3-oxohexanoyl)-L-homoserine lactone, and pathogenicity in
            soft-rotting Erwinia spp
  JOURNAL  J. Bacteriol. 177(17) (1995) In press
COMMENT    Original source text: Erwinia carotovora (strain 71, sub_species
            carotovora) DNA.
FEATURES   Location/Qualifiers
  source   1..500
            /organism="Pectobacterium carotovorum"
            /strain="71"
            /sub_species="carotovora"
            /db_xref="taxon:554"
  gene     107..431
            /gene="rsmA"
  -10_signal 107..112
            /gene="rsmA"
  RBS      235..239
            /gene="rsmA"
  CDS      246..431
            /gene="rsmA"
            /function="global repressor"
            /note="putative"
            /codon_start=1
            /transl_table=11
            /protein_id="AAA74502.1"
            /db_xref="GI:927032"
            /translation="MLIILTRRVGETLIIGDEVTVTVLGVKGNQVRIGVNAPKEVSVHR
            EEIYQRIQAEKSQPTSY"
BASE COUNT 140 a    101 c    120 g    139 t
ORIGIN
  1  ggatccggca agcaggatag aaagtgtggtt accttcagat attctgaagc tttacatgct
  61  cagttctggt gttgtgataa caaaagcaca agctactgat atcgactaaa ctaacaagta
 121  gtgacaaacc ggagtgtgat ggtgtgggta taccatcgtc taggtttacg ttttcacagc
 181  acatgatgga taatggcggg gagacagaga gacccgactc tttataatct ttcaaggagc
 241  aaagaatgct tattttgact cgtcgagttg gcgaaacct catcatcggc gatgaggtaa
 301  cggttaccgt attaggagtg aaaggcaacc aggtgcgtat tgggtgtaat gcacctaaag
 361  aggtttctgt ccaccgtgaa gagatctatc agcgtattca ggccgaaaaa tctcaaccaa
 421  cgtcatattg attgacaatg cgtctcgtgt tcgcgggacg caattgttat ttccggtttt
 481  tccccacac atttatcgat
//

```

The two entries shown above contain the same biological information but differ in the format and presentation of this information. One of the most obvious difference is in the header region of the file that gives the background information to the submitted sequence. Another clear difference is that the EMBL entry has an additional **two letter line code** on the left hand margin.

EMBL entries are structured so as to be usable by human readers as well as by computer programs. The explanations, descriptions, classifications and other comments are in ordinary English for readability. At the same time, the structure is systematic enough to allow computer programs to easily read, identify, and manipulate the various types of data included.

Each line begins with a **two letter line code**, which indicates the type of information contained in the line. The currently used line types, along with their respective line codes, are listed below.

| | |
|--------------------------------|----------------------------------|
| ID - identification | (begins each entry; 1 per entry) |
| AC - accession number | (>=1 per entry) |
| SV - new sequence identifier | (>=1 per entry) |
| DT - date | (2 per entry) |
| DE - description | (>=1 per entry) |
| KW - keyword | (>=1 per entry) |
| OS - organism species | (>=1 per entry) |
| OC - organism classification | (>=1 per entry) |
| OG - organelle | (0 or 1 per entry) |
| RN - reference number | (>=1 per entry) |
| RC - reference comment | (>=0 per entry) |
| RP - reference positions | (>=1 per entry) |
| RX - reference cross-reference | (>=0 per entry) |
| RA - reference author(s) | (>=1 per entry) |
| RT - reference title | (>=1 per entry) |
| RL - reference location | (>=1 per entry) |
| DR - database cross-reference | (>=0 per entry) |
| FH - feature table header | (0 or 2 per entry) |
| FT - feature table data | (>=0 per entry) |
| CC - comments or notes | (>=0 per entry) |
| XX - spacer line | (many per entry) |
| SQ - sequence header | (1 per entry) |
| bb - (blanks) sequence data | (>=1 per entry) |
| // - termination line | (ends each entry; 1 per entry) |