

# References

Abbot, J. C. et al. (2005) *Bioinformatics* 21(18) 3665-3666. WebACT – an online companion for the Artemis Comparison Tool.

Allen JE & Salzberg SL (2005). *Bioinformatics* 21: 3596-3603. JIGSAW: integration of multiple sources of evidence for gene prediction.

Alexa A. et al. (2006) *Bioinformatics* 22: 1600-1607. Improved scoring of functional groups from gene expression data by decorrelating GO graph structure.

Anders S & Huber W (2010) *Genome Biol* 11: R106. Differential expression analysis for sequence count data.

Anders S. HTSeq: Analysing high-throughput sequencing data with Python. 2010. Software. [<http://www-huber.embl.de/users/anders/HTSeq/>]

Assefa, S. et al. (2009) *Bioinformatics* 25 (15) 1968-9. ABACAS: algorithm-based automatic contiguation of assembled sequences.

Berriman, M., and K. Rutherford (2003) *Brief Bioinform* 4 (2) 124-132. Viewing and annotating sequence data with Artemis.

Bozdech Z. et al. (2003) *PLOS Biol* 1: E5. The transcriptome of the intraerythrocytic developmental cycle of *Plasmodium falciparum*.

Carver T. J. et al. (2010) *Bioinformatics* (doi:10.1093/bioinformatics/btq010)  
BamView: Viewing mapped read alignment data in the context of the reference sequence.

Carver T. J. et al. (2005) *Bioinformatics* 21: 3422-3. ACT: the Artemis Comparison Tool.

Conesa A, et al. (2005) *Bioinformatics* 21: 3674-3676. Blast2GO: a universal tool for annotation, visualization and analysis in functional genomics research.

- Delcher AL. et al. (1999) *Nucleic Acids Res* 27: 4636-4641. Improved microbial gene identification with GLIMMER.
- Gardner et al. (2002). *Nature* 419(6906):498-511. Genome sequence of the human malaria parasite *Plasmodium falciparum*.
- Grant GR, et al. (2011) *Bioinformatics* 27: 2518-2528. Comparative analysis of RNA-Seq alignment algorithms and the RNA-Seq unified mapper (RUM).
- Hacker, J. et al. (1997) *Mol Microbiol* 23: 1089-97. Pathogenicity islands of virulent bacteria: structure, function and impact on microbial evolution.
- Haas BJ, et al. (2008). *Genome Biol* 9: R7. Automated eukaryotic gene structure annotation using EVIDENCEModeler and the Program to Assemble Spliced Alignments.
- Hardcastle TJ & Kelly KA (2010). *BMC Bioinformatics* 11: 422. baySeq: empirical Bayesian methods for identifying differential expression in sequence count data
- Kozarewa, I., Z. Ning, et al. (2009). *Nature Met* 6(4): 291-295. Amplification-free Illumina sequencing-library preparation facilitates improved mapping and assembly of (G+C)-biased genomes.
- Langmead et al. (2009). *Genome Biol* 10:R25. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome.
- Li et al. (2008a). *Genome Res* 18:1851-8. Mapping short DNA sequencing reads and calling variants using mapping quality scores.
- Li et al. (2008b). *Bioinformatics* 24(5):713-714. SOAP: short oligonucleotide alignment program.
- Li et al. (2009). *Bioinformatics*, 25:1754-60. Fast and accurate short read alignment with Burrows-Wheeler Transform.
- Majoros et al. (2003) *Nucleic Acids Res* 31 (13) 3601-3604. GlimmerM, Exonomy and Unveil: three *ab initio* eukaryotic genefinders.
- Mortazavi et al. (2008). *Nature Met* 5: 621 – 628. Mapping and quantifying mammalian transcriptomes by RNA-Seq.
- Ning et al. (2001). *Genome Res* 10:1725-9. SSAHA: a fast search method for large DNA databases.

- 
- Otto et al. (2010) *Mol Microbiol* Apr;76(1):12-24. New insights into the blood stage transcriptome of *Plasmodium falciparum* using RNA-Seq.
- Otto, T. D., G. P. Dillon, et al. (2011). *Nucleic Acids Res* **39**(9): e57. RATT: Rapid Annotation Transfer Tool.
- Otto, T. D., M. Sanders, et al. (2010). *Bioinformatics* **26**(14): 1704-1707. Iterative Correction of Reference Nucleotides (iCORN) using second generation sequencing technology.
- Parkhill, J. (2002) *Method Microbiol* 33: 1-26. Annotation of microbial genomes.
- Robinson MD, et al. (2010) *Bioinformatics* 26: 139-140. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data.
- Rutherford et al. (2000) *Bioinformatics* 16 (10) 944-945. Artemis: sequence visualization and annotation.
- Simpson, J. T., K. Wong, et al. (2009). *Genome Res* **19**(6): 1117-1123. ABySS: a parallel assembler for short read sequence data.
- Stephens et al. (1998). *Science* 282(5389): 754 – 759. Genome sequence of an obligate intracellular pathogen of humans: *Chlamydia trachomatis*.
- Tsai, I. J., T. D. Otto, et al. (2010). *Genome Biol* **11**(4): R41. Improving draft assemblies by iterative mapping and assembly of short reads to eliminate gaps.
- Trapnell et al. (2009). *Bioinformatics* 25(9):1105-1111. TopHat: discovering splice junctions with RNA-Seq.
- Wang et al. (2009). *Nat Rev Genet* 10(1):57-63. RNA-Seq: A revolutionary tool for transcriptomics.
- Zerbino, D. R. and E. Birney (2008). *Genome Res* **18**(5): 821-829. Velvet: algorithms for *de novo* short read assembly using de Bruijn graphs.

# Appendices

**Appendix I:** Course Virtual Machine (VM) Quick Start Guide

**Appendix III:** ACT comparison files

**Appendix IV:** Feature Keys and Qualifiers – a brief explanation of what they are and a sample of the ones we use.

**Appendix V:** Generating ACT comparison files using BLAST

**Appendix VIII:** Prokaryotic Protein Classification Scheme used within the PSU

**Appendix IX:** List of colour codes

**Appendix XI:** Splice site information

## Appendix I: Course Virtual Machine (VM) Quick Start Guide

Using a VM enables us to encapsulate the course data and software in such a way that you can still make use of them when you return to your own laboratory.

To use the VM on the USB stick provided, you will first need to download VirtualBox (<http://www.virtualbox.org/>). This software is required to run the VM on your machine, it is free and available for windows, MacOSX and linux,

For a detailed description of VirtualBox and the installation see the on-line manual (<http://www.virtualbox.org/manual/>).

### Download and Install VirtualBox

- Download VirtualBox for the type of workstation you are using (e.g. Windows) from <http://www.virtualbox.org/wiki/Downloads>.
- Double click on the executable file (Windows). The installation welcome dialog opens and allows you to choose where to install VirtualBox to, and which components to install. Depending on your Windows configuration, you may see warnings about "unsigned drivers" or similar. Please select "Continue" on these warnings; otherwise VirtualBox might not function correctly after installation.
- Launch the VirtualBox software from the desktop shortcut or from the program menu.

### Setting up the VM

VirtualBox needs to be pointed at the VDI (This is the file that is on the memory stick used during the course) file as follows:

- Insert the USB memory stick provided. This contains a Virtual Disk Image (VDI) file.

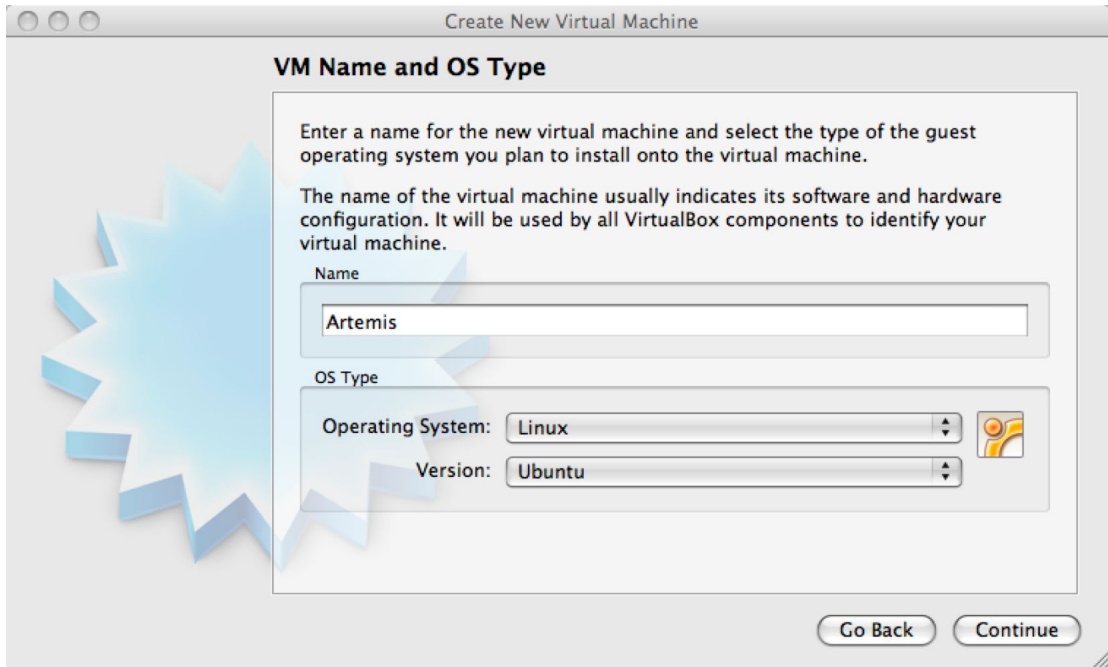
Create a new virtual machine by selecting 'New' from the options at the top. Then fill the boxes in as shown below:

In the first window enter:

Name: **Artemis**

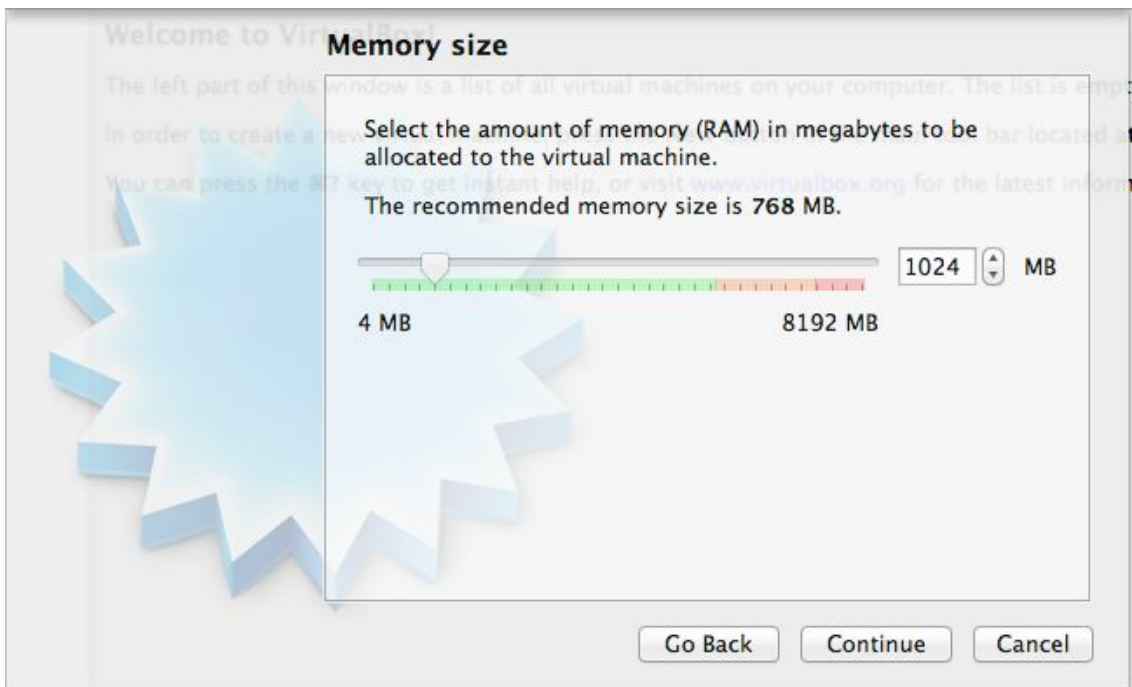
Operating System: **Linux**

Version: **Ubuntu**



Click 'Continue'

In the next window set the memory to at least 1GB (as shown), but 2GB (2048 MB) will give you better performance. You can use more but no more than half the amount of memory on your PC.



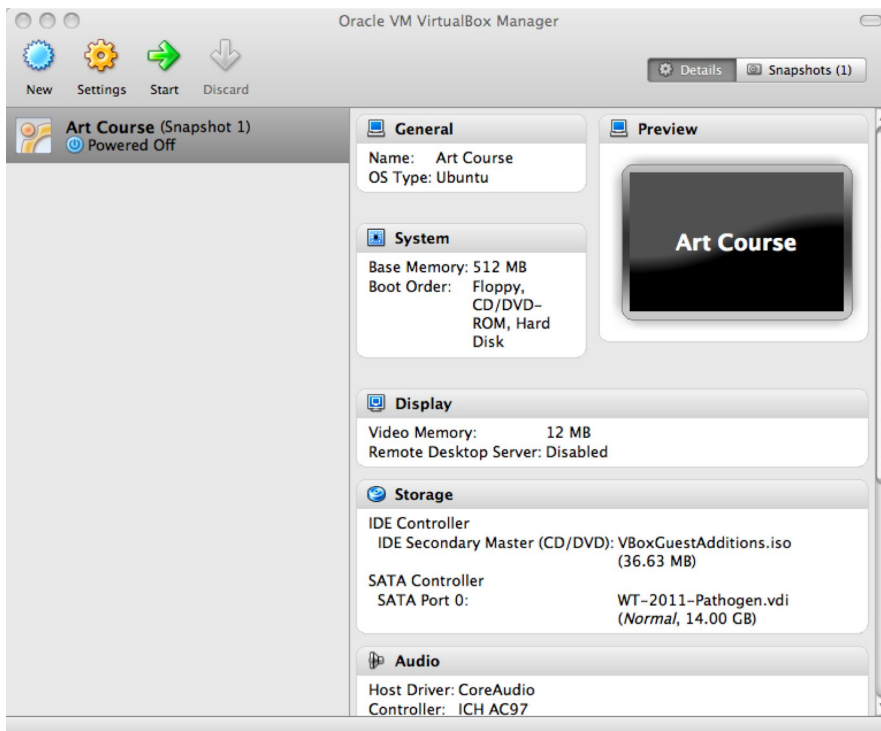
Click 'Continue'.

In the next window select ‘Use existing hard disk’ and from the folder icon on the right hand side navigate to the memory USB stick and select the VDI file located on the memory stick



Click ‘Continue’.

There will now be an ‘Artemis’ (powered off) button in the left hand side of VirtualBox.



Double click on this new Artemis course power button to start the VM. It will then log you into the Ubuntu desktop.

### **Setting up a Shared Folder**

This allows you to share a folder between the VM and your workstation. This means you can put files that you want to share between the operating systems in this folder.

Create a directory to share called 'VMshare' on your machine. With the VM shutdown select the 'Artemis' button in VirtualBox and click 'Settings' in the top menu bar. Go to 'Shared Folders' and select the '+' button on the right. In the 'Folder Path' select 'Other' and navigate to and select the 'VMshare' folder that you have created. Then click on 'OK'.

When the 'Artemis' VM is next started double click on the 'mount' icon in your home folder. This will open a window that you need to type the password into:

```
wt
```

It will show the contents of this folder in the /home/wt/host directory in Ubuntu.

### **A note on memory usage:**

Some computing processes are very memory hungry. Should you find that your computer processes are killed without a clear reason, one aspect to check is the amount of memory allocated to the VM. The 1024MB you have allocated using this tutorial has been check and should be enough. Nonetheless, the amount of memory allocated to the VM can be changed at any time.



### Appendix III: ACT comparison files

ACT supports three different comparison file formats:

- 1) BLAST version 2.2.2 output: The blastall command must be run with the -m 8 flag which generates one line of information per HSP.
- 2) MegaBLAST output: ACT can also read the output of MegaBLAST, which is part of the NCBI blast distribution.
- 3) MSPcrunch output: MSPcrunch is program for UNIX and GNU/Linux systems which can post-process BLAST version 1 output into an easier to read format. ACT can only read MSPcrunch output with the -d flag.

Here is an example of an ACT readable comparison file generated by MSPcrunch -d.

```
1399 97.00 940 2539 sequence1.dna 1 1596 AF140550.seq
1033 93.00 9041 10501 sequence1.dna 9420 10880 AF140550.seq
828 95.00 6823 7890 sequence1.dna 7211 8276 AF140550.seq
773 94.00 2837 3841 sequence1.dna 2338 3342 AF140550.seq
```

The columns have the following meanings (in order): score, percent identity, match start in the query sequence, match end in the query sequence, query sequence name, subject sequence start, subject sequence end, subject sequence name.

The columns should be separated by single spaces.

## Appendix IV: Feature Keys and Qualifiers – a brief explanation of what they are and a sample of the ones we use.

**1 – Feature Keys:** They describe features with DNA coordinates and once marked, they all appear in the Artemis main window. The ones we use are:

**CDS:** Marks the extent of the coding sequence.

**RBS:** Ribosomal binding site

**misc\_feature:** Miscellaneous feature in the DNA

**rRNA:** Ribosomal RNA

**repeat\_region**

**repeat\_unit**

**stem\_loop**

**tRNA:** Transfer RNA

**2 – Qualifiers:** They describe features in relation to their coordinates. Once marked they appear in the lower part of the Artemis window. They describe the feature whose coordinates appear in the ‘location’ part of the editing window. The ones we commonly use for annotation at the Sanger Institute are:

**/class:** Classification scheme we use “in-house” developed from Monica Riley’s MultiFun assignments (see Appendix VI).

**/colour:** Also used in-house in order to differentiate between different types of genes and other features.

**/gene:** Descriptive gene name, eg. *ilvE*, *argA* etc.

**/label:** Allows you to label a gene/feature in the main view panel.

**/note:** This qualifier allows for the inclusion of free text. This could be a description of the evidence supporting the functional prediction or other notable features/information which cannot be described using other qualifiers.

**/product:** The assigned possible function for the protein goes here.

**/pseudo:** Matches in different frames to consecutive segments of the same protein in the databases can be linked or joined as one and edited in one window. They are marked as pseudogenes. They are normally not functional and are considered to have been mutated.

**/locus\_tag :** Systematic gene number, eg SAS1670, Sty2412 etc.

The list of keys and qualifiers accepted by EMBL in sequence/annotation submission files are list at the following web page:

<http://www3.ebi.ac.uk/Services/WebFeat/>

## Appendix V: Generating ACT comparison files using BLAST

The following pages demonstrate how you can generate your own comparison files for ACT from a stand-alone version of the BLAST software. In Appendix X the NCBI BLAST distribution was downloading onto a PC with Windows XP. The exercises in this module are based on the Linux version of the BLAST software. Although the operating systems are different, the command lines used to run the programs are the same. One of the main differences between the two operating systems is that in Windows the BLAST program command line is run in the DOS Command Prompt window, whereas in Linux it is run from a Xterminal window.

In the exercises below you are going to download two small sequences (plasmids), and for two large sequences (whole genomes). You are then going generate files containing DNA sequences in FASTA format for these sequences, which will then be compared using two different programs from the NCBI BLAST distribution to generate ACT comparison files.

### Exercise 1

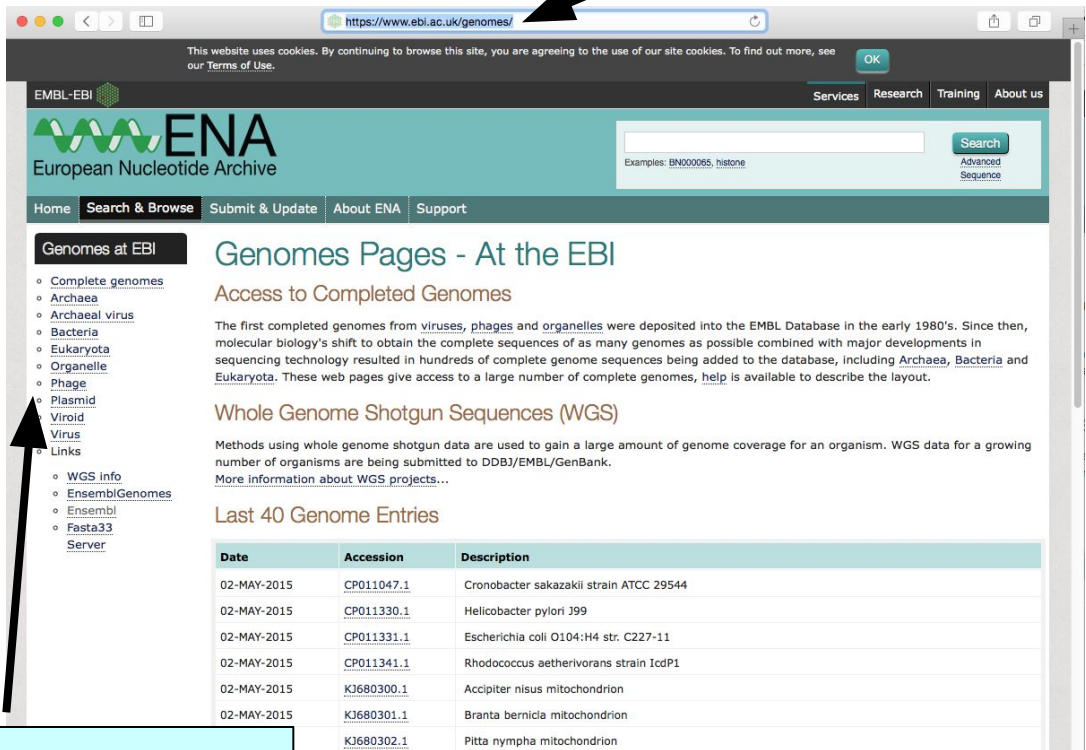
In this exercise you are going to download two plasmid sequences in EMBL format from the EBI genomes web page. You are then going to use Artemis to write out the DNA sequences of both plasmids in FASTA format. These two FASTA format sequences will then be compared using the blastall program from the NCBI BLAST distribution. Using blastall you can run BLASTN to identify regions of DNA-DNA similarity and write out a ACT readable comparison file. If required, blastall can also be used to run other flavours of BLAST with the appropriate input files (i.e. DNA files for TBLASTX, protein files for BLASTP, and protein and DNA for BLASTX). For the purpose of generating ACT comparison files BLASTN and TBLASTX are appropriate.

In this example two relative small sequences have been chosen (<500 kb). BLAST running on a relatively modern stand alone machine can easily deal with required computations, and thus the comparison file should be produced in a matter of seconds. However as the size of the compared sequences increases the time taken to produce the output will dramatically increase. Therefore for very large sequences (several Mb) it will be impractical to run them using blastall. In **Exercise 2** you will use megablast, another program in the NCBI BLAST distribution, which is useful for comparing large sequence that are very similar.

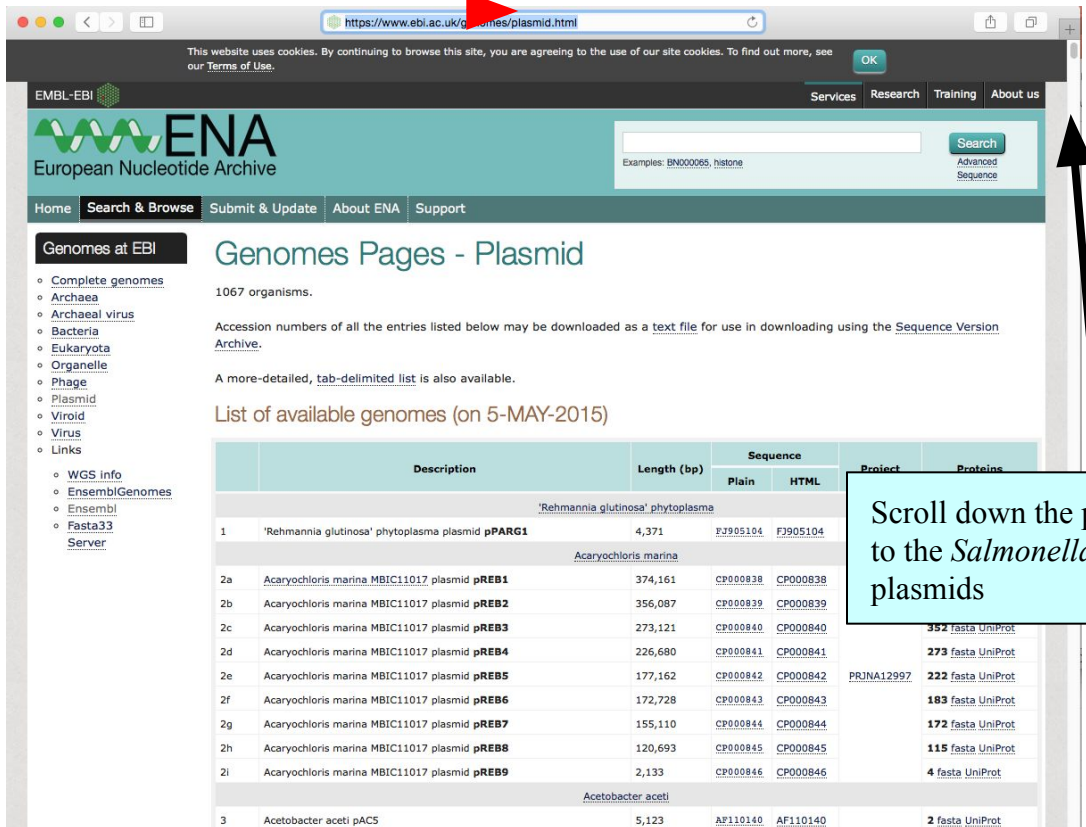
The plasmids chosen for this comparison are the multiple drug resistance incH1 plasmid pHCM1 from the sequenced strain of *Salmonella typhi* CT18 originally isolated in 1993, and R27, another incH1 plasmid first isolated from *S. typhi* in the 1960s.

# Downloading the *S. typhi* plasmid sequences

Go to the EBI genomes web page (<http://www.ebi.ac.uk/genomes>)



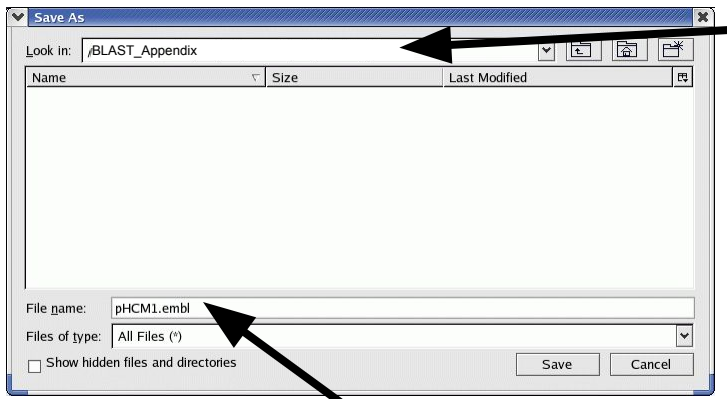
Click on the Plasmid hyperlink



Scroll down the page to the *Salmonella* plasmids

Press the Shift key and left Click on the accession number hyperlink for pHCM1 (AL513383) in the Plain Sequence column

Accession	Plasmid Name	Size (kb)	FASTA SRS
158a	Riemerella anatipestifer plasmid pCFC1	3,966	4 FASTA SRS
158b	Riemerella anatipestifer plasmid pCFC2	5,609	3 FASTA SRS
159	Ruminococcus flavefaciens R13e2 cryptic plasmid pBAW301	1,768	1 FASTA SRS
160	Salmonella choleraesuis strain 79500 plasmid pSFD10	4,871	6 FASTA SRS
161	Salmonella enterica subsp. enterica serovar Berta plasmid pBERT	4,656	9 FASTA SRS
162a	Salmonella enterica subsp. enterica serovar Typhi str. CT18 plasmid pHCM1	218,160	234 FASTA SRS
162b	Salmonella enterica subsp. enterica serovar Typhi str. CT18 plasmid pHCM2	106,516	132 FASTA SRS
163	Salmonella enterica subsp. enterica serovar Typhimurium plasmid pFPTB1	12,656	6 FASTA SRS
164a	Salmonella enteritidis serovar Enteritidis plasmid pC	5,269	4 FASTA SRS
164b	Salmonella enteritidis serovar Enteritidis plasmid pK	4,245	3 FASTA SRS
164c	Salmonella enteritidis serovar Enteritidis plasmid pP	4,301	3 FASTA SRS
165a	Salmonella typhi R27 plasmid	180,461	204 FASTA SRS
165b	Salmonella typhi plasmid R27	38,245	34 FASTA SRS



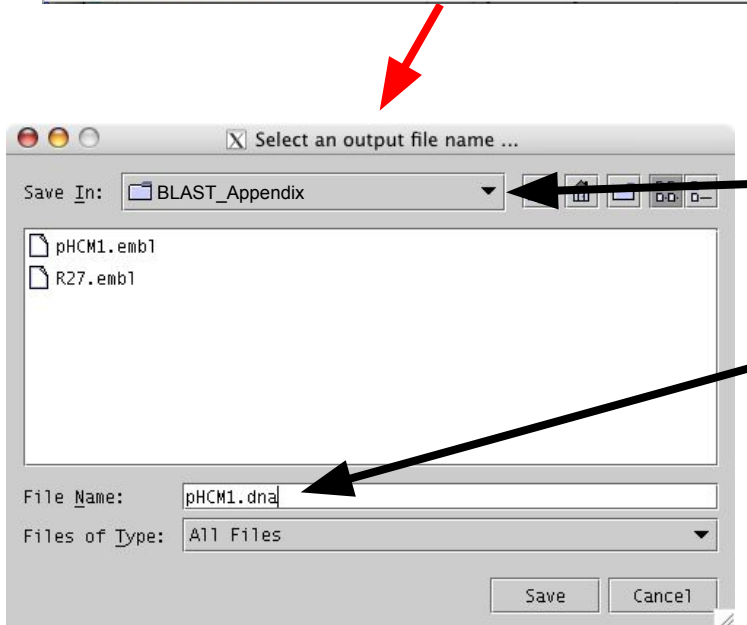
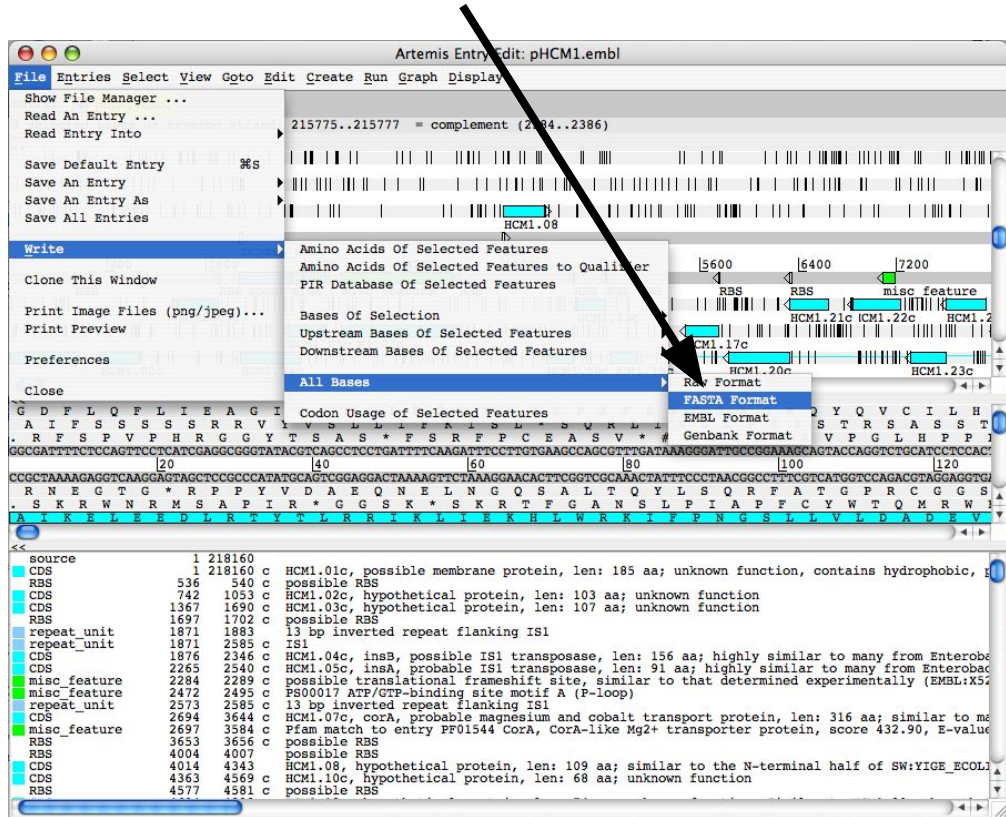
Save the EMBL sequence in a suitable directory. For example: BLAST\_Appendix

Save the file as pHCM1.embl

Repeat for the *Salmonella typhi* R27 plasmid (AF250878). Be careful when choosing the plasmid to download as there is also a *Salmonella typhi* plasmid R27 entry (AF105019), the one that you want is the larger of the two, 180,461 kb as opposed to 38,245 kb – make sure the accession number is correct. Save as R27.embl.

In order to run BLASTN you require two DNA sequences in FASTA format. The pHCM1 and R27 sequences previously downloaded from the EBI are EMBL format files, i.e. they contain protein coding information and the DNA sequence. In order to generate the DNA files in FASTA format, Artemis can be used as follows.

Load up the plasmid EMBL files in **Artemis** (each plasmid requires a separate Artemis window), select **Write, All Bases, FASTA format**.



Save the DNA sequence in the BLAST\_Appendix directory

Save as pHCM1.dna

Also do this for R27.embl

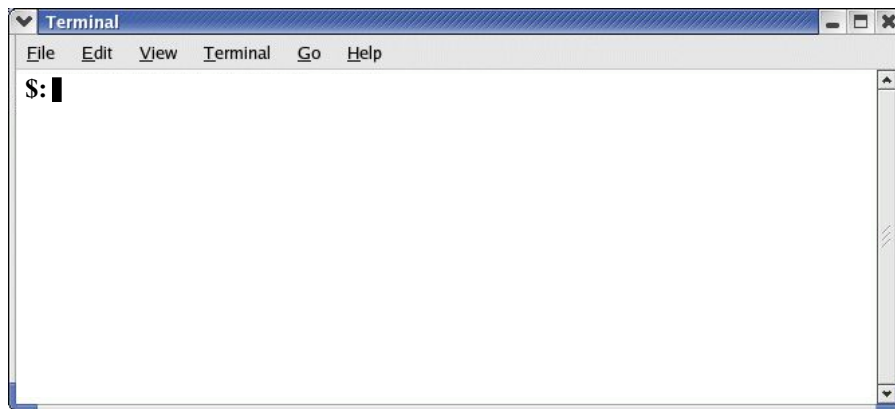
## Running Blast

There are several programs in the BLAST package that can be used for generating sequence comparison files. For a detailed description of the uses and options see the appropriate README file in the BLAST software directory (see Appendix X).

In order to generate comparison files that can be read into ACT you can use the **blastall** program running either BLASTN (DNA-DNA comparison) or TBLASTX (translated DNA-translated DNA comparison) protocols.

As an example you will run a BLASTN comparison on two relatively small sequences; the pHCM1 and R27 plasmids from *S. typhi*. In principle any DNA sequences in FASTA format can be used, although size becomes an issue when dealing with sequences such as whole genomes of several Mb (see **Exercise 2** in this module). When obtaining nucleotide sequences from databases such as EMBL using a server such as SRS (<http://srs.ebi.ac.uk>), it is possible to specify that the sequences are in FASTA format.

To run the BLAST software you will need an Xterminal window like the one below. If you do not already have one opened, you can open a new window by clicking on the Xterminal icon on the menu bar at the bottom of your screen.



Make sure you are in the appropriate directory (in this example it is BLAST\_Appendix.) You should now see both the new FASTA files for the pHCM1 and R27 sequences in the BLAST\_Appendix directory as well as their respective EMBL format files. (Hint: You can use the **pwd** command to check the present working directory, the **cd** command to change directories, and the **ls** command will list the contents of the present working directory).

When comparing sequences in BLAST, one sequence is designated as a **database** sequence, and the other the **query** sequence. Before you run BLAST you have to format one of the sequences so that BLAST recognises it as a database sequence. **formatdb** is a program that does this and comes as part of the NCBI BLAST distribution.

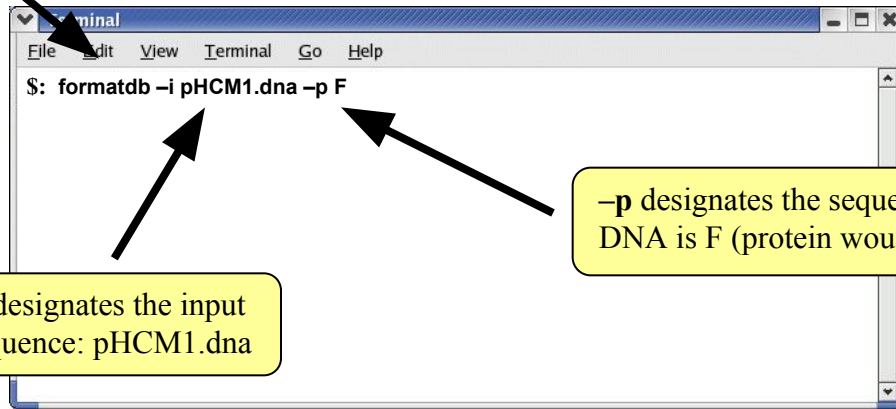
You will treat pHCM1.dna as the **database** sequence and R27.dna as the **query** sequence

At the Command Prompt type:  
**formatdb -i pHCM1.dna -p F**  
Press **Return**

**formatdb** is the database format program

**-i** designates the input sequence: pHCM1.dna

**-p** designates the sequence type: DNA is F (protein would be T)



Now you can run the BLAST on the two plasmid sequences. The program that you are going to use is **blastall**. In addition to the standard command line inputs we have to add an additional flag (**-m 8**) to the command line so that the BLAST output can be read by ACT. This specifies that the output of BLAST is in one line per entry format (see appendix II).

At the Command Prompt type:  
**blastall -p blastn -m 8 -d pHCM1.dna -i R27.dna -o pHCM1\_vs\_R27**  
Press **Return**

**tblastx** could be substituted here if a translated DNA-translated DNA comparison was required

**-o** designates the output file: pHCM1\_vs\_R27

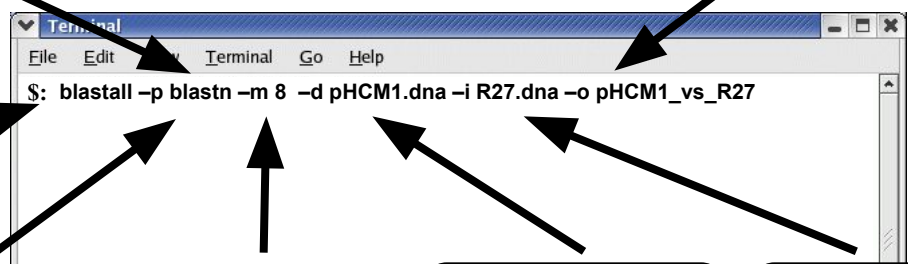
**blastall** is the BLAST program

**-p** designates the flavour of BLAST: **blastn** (in this instance a DNA-DNA comparison)

**-m 8** designates the ACT readable output

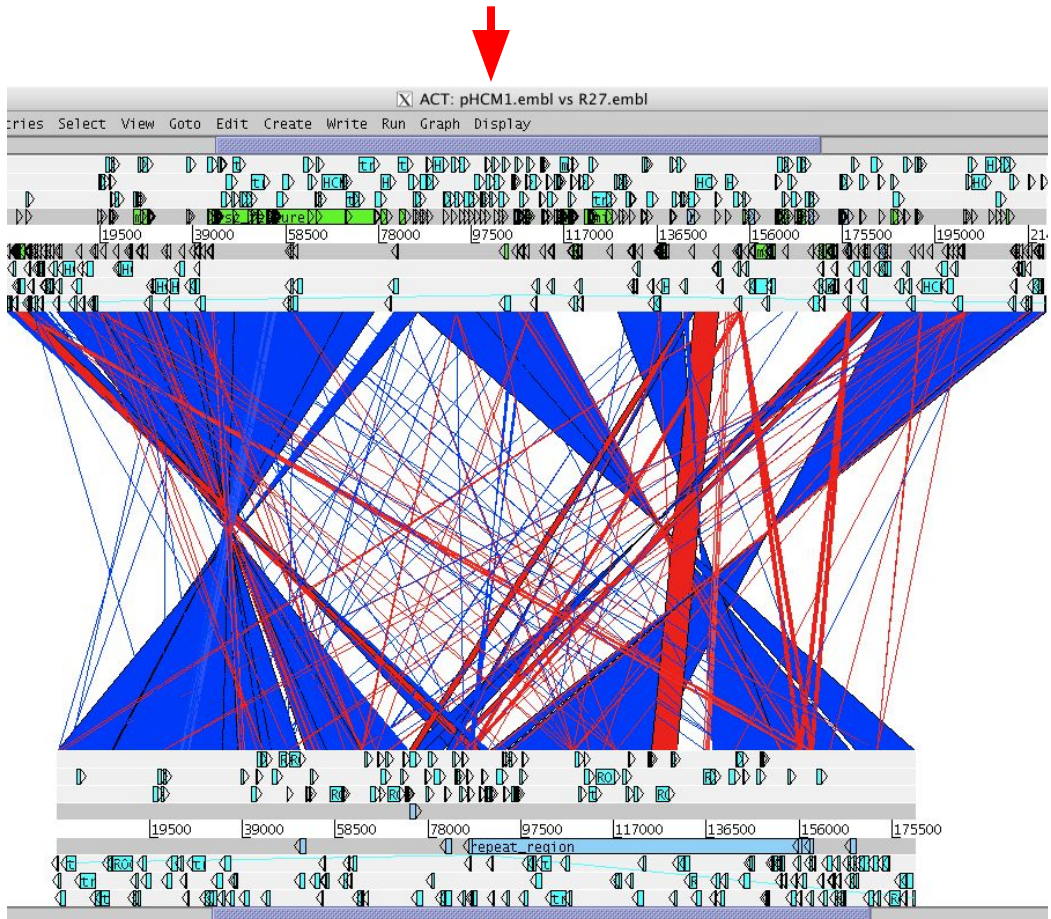
**-d** designates the database sequence: pHCM1.dna

**-i** designates the query sequence: R27.dna





The pHCM1\_vs\_R27 comparison file can now be read into ACT along with the pHCM1.embl and R27.embl (or pHCM1.dna and R27.dna) sequence files.



The result of the BLASTN comparison shows that there are regions of DNA shared between the plasmids; pHCM1 shares 169 kb of DNA at greater than 99% sequence identity with R27. Much of the additional DNA in the pHCM1 plasmid appears to have been inserted relative to R27 and encodes functions associated with drug resistance. What antibiotic resistance genes can you find in the pHCM1 plasmid that are not found in R27?

The two plasmids were isolated more than 20 years apart. The comparison suggests that there have been several independent acquisition events that are responsible for the multiple drug resistance seen in the more modern *S. typhi* plasmid.

## Exercise 2

In the previous exercise you used BLASTN to generate a comparison file for two relatively small sequences (>500,000 kb). In the next exercise we are going to use another program from NCBI BLAST distribution, **megablast**, that can be used for nucleotide sequence alignment searches, i.e. DNA-DNA comparisons. If you are comparing large sequences such as whole genomes of several Mb, the **blastall** program is not suitable. The BLAST algorithms will struggle with large DNA sequences and therefore the processing time to generate a comparison file will increase dramatically.

**megablast** uses a different algorithm to BLAST which is not as stringent which therefore makes the program faster. This means that it is possible to generate comparison files for genome sequences in a matter of seconds rather than minutes and hours.

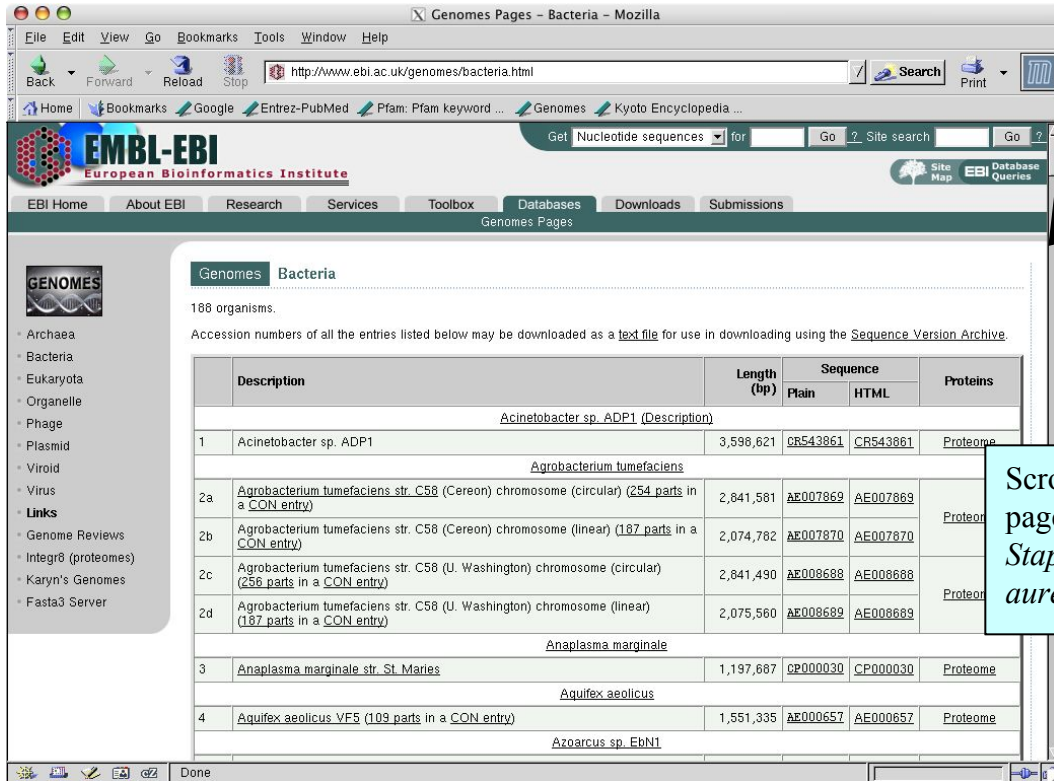
There are some drawbacks to using this program. Firstly, only DNA-DNA alignments (BLASTN) can be performed using **megablast**, rather than translated DNA-DNA alignments (TBLASTX) as can be using **blastall**. Secondly as the algorithm used is not as stringent, **megablast** is suited to comparing sequences with high levels of similarity such as genomes from the same or very closely related species.

In this exercise you are going to download two *Staphylococcus aureus* genome sequences from the EBI genomes web page and use Artemis to write out the FASTA format DNA sequences for both as before in **Exercise 1**. These two FASTA format sequences will then be compared using **megablast** to identify regions of DNA-DNA similarity and write out an ACT readable comparison file.

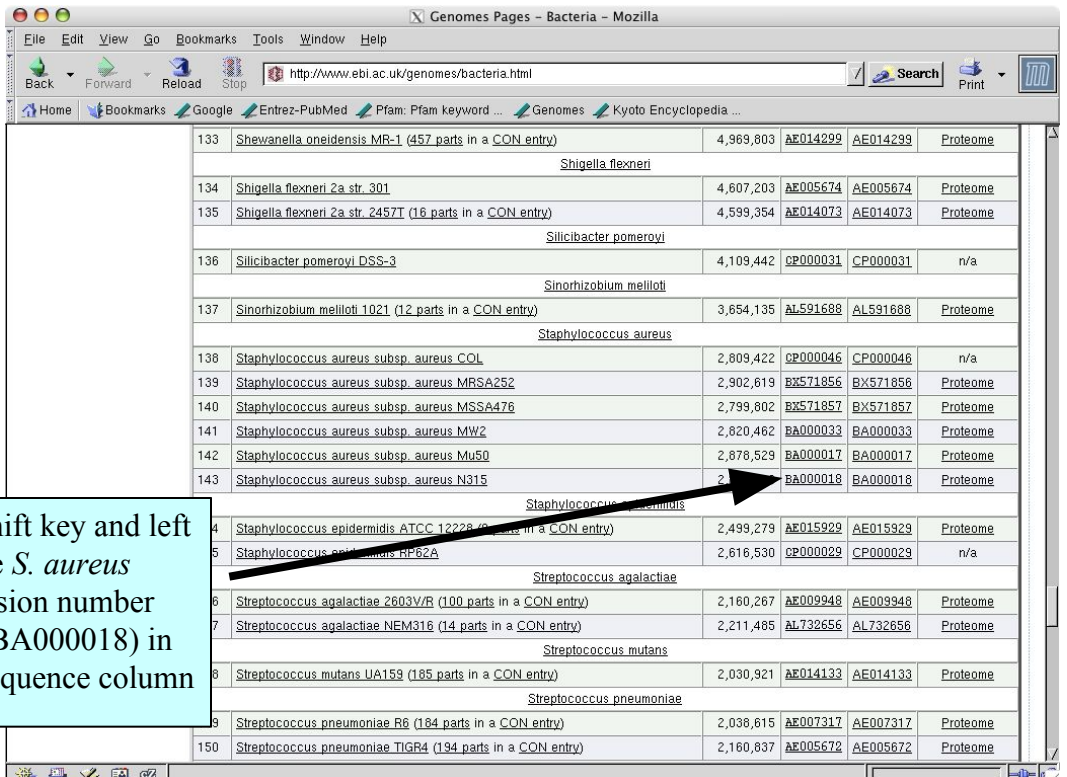
The genomes that have been chosen for this comparison are from a hospital-acquired methicillin resistant *S. aureus* (MRSA) strain N315 (BA000018), and a community-acquired MRSA strain MW2 (BA000033).

## Downloading the *S. aureus* genomic sequences

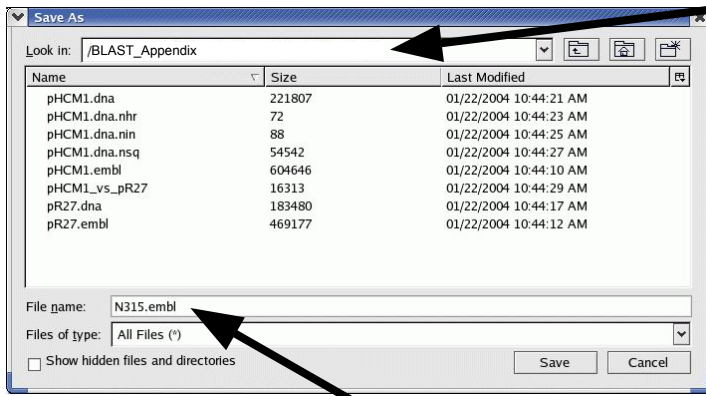
Go to the EBI genomes web page (<http://www.ebi.ac.uk/genomes>) as before in Exercise 2, and click on the **Bacteria** hyperlink



Scroll down the page to the *Staphylococcus aureus* genomes



Press the Shift key and left Click on the *S. aureus* N315 accession number hyperlink (BA000018) in the Plain Sequence column



Save the EMBL sequence  
in a suitable directory.  
For example:  
BLAST\_Appendix

Save the file as N315.embl

Repeat for the *S. aureus* MW2 genome (BA000033). Be careful when choosing the genome to download as there is another *S. aureus* genome entry for strain Mu50 (BA000017). Save as MW2.embl.

Generate DNA files in FASTA format using Artemis for both the genome sequences as previously done in exercise 1.

(Hint: In **Artemis** (each genome requires a separate Artemis window), select **Write, Write All Bases, FASTA format**).

Save the DNA sequences as N315.dna and MW2.dna for the respective genomes.

## Running Blast

In the previous exercise you used the **blastall** program to run BLASTN on two plasmid sequences. As the genome sequences are larger (~2.8 Mb) you are going to run **megablast**, another program from the NCBI BLAST distribution that can generate comparison files in a format that ACT can read (see Appendix II). For a detailed description of the uses and options in **megablast** see the megablast README file in the BLAST software directory (Appendix X).

As before you will run the program from the command line in an Xterminal window.

Like BLAST, **megablast** requires that one sequence is designated as a **database** sequence and the other the **query** sequence. Therefore one of the sequences has to be formatted so that Blast recognises it as a database sequence. This can be done as before using **formatdb**.

We will treat N315.dna as the **database** sequence and MW2.dna as the **query** sequence

At the Command Prompt type:  
**formatdb -i N315.dna -p F**  
Press **Return**

```
Terminal
File Edit View Terminal Go Help
$: formatdb -i N315.dna -p F
```

**-i** designates the input sequence: N315.dna

**-p** designates the sequence type: DNA is F (protein would be T)

Now we can run the **megablast** on the two MRSA genome sequences. The default output format is one line per entry that ACT can read, therefore there is no need to add an additional flag (i.e. -m 8) to the command line (see appendix II).

At the Command Prompt type:  
**megablast -d N315.dna -i MW2.dna -o N315\_vs\_MW2**  
Press **Return**

**megablast** is the program

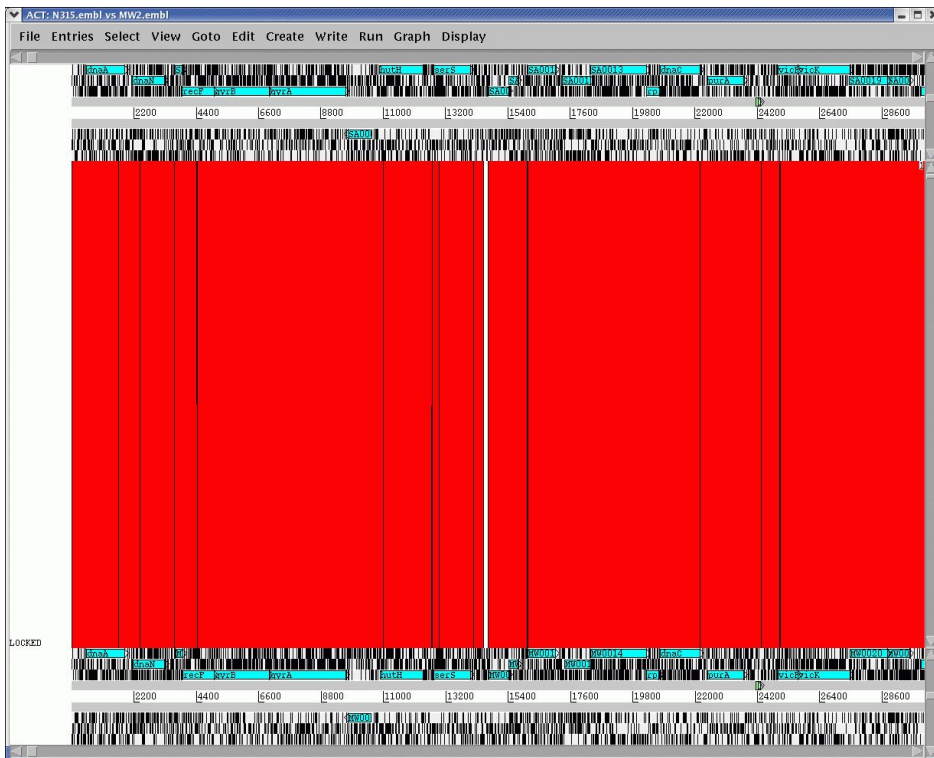
**-o** designates the output file: N315\_vs\_MW2

```
Terminal
File Edit View Terminal Go Help
$: megablast -d N315.dna -i MW2.dna -o N315_vs_MW2
```

**-d** designates the database sequence: N315.dna

**-i** designates the query sequence: MW2.dna

The N315\_vs\_MW2 comparison file can now be read into ACT along with the N315.embl and MW2.embl (or N315.dna and MW2.dna) sequence files.



A comparison of the N315 and MW2 genomes in ACT using the **megablast** comparison reveals a high level of synteny (conserved gene order). This is perhaps not unsurprising as both genomes belong to strains of the same species. Using results of comparisons like these it is possible to identify genomic differences that may contribute to the biology of the bacteria and also investigate mechanisms of evolution.

Both N315 and MW2 are MRSA, however N315 is associated with disease in hospitals, and MW2 causes disease in the community and is more invasive. Scroll rightward in both genomes to find the first large region of difference. Examine the annotation for the genes in these regions. What are the encoded functions associated with these regions? What significance does this have for the evolution of methicillin resistance in these two *S. aureus* strains from clinically distinct origins?

## Appendix VIII: Prokaryotic Protein Classification Scheme used within the PSU

This scheme was adapted for Sanger in-house use from the Monica Rileys protein classification (<http://genprotec.mbl.edu/files/Multifun.html>).

More classes can be added depending on the microorganism that is being annotated (e.g secondary metabolites, sigma factors (ECF or non-ECF), etc).

- 0.0.0 Unknown function, no known homologs
- 0.0.1 Conserved in Escherichia coli
- 0.0.2 Conserved in organism other than Escherichia coli
- 1.0.0 Cell processes
  - 1.1.1 Chemotaxis and mobility
  - 1.2.1 Chromosome replication
  - 1.3.1 Chaperones
  - 1.4.0 Protection responses
    - 1.4.1 Cell killing
    - 1.4.2 Detoxification
    - 1.4.3 Drug/analog sensitivity
    - 1.4.4 Radiation sensitivity
  - 1.5.0 Transport/binding proteins
    - 1.5.1 Amino acids and amines
    - 1.5.2 Cations
    - 1.5.3 Carbohydrates, organic acids and alcohols
    - 1.5.4 Anions
    - 1.5.5 Other
  - 1.6.0 Adaptation
    - 1.6.1 Adaptations, atypical conditions
    - 1.6.2 Osmotic adaptation
    - 1.6.3 Fe storage
  - 1.7.1 Cell division
- 2.0.0 Macromolecule metabolism
- 2.1.0 Macromolecule degradation
  - 2.1.1 Degradation of DNA
  - 2.1.2 Degradation of RNA
  - 2.1.3 Degradation of polysaccharides
  - 2.1.4 Degradation of proteins, peptides, glycoproteins
- 2.2.0 Macromolecule synthesis, modification
  - 2.2.01 Amino acyl tRNA synthesis; tRNA modification
  - 2.2.02 Basic proteins - synthesis, modification
  - 2.2.03 DNA - replication, repair, restriction./modification
  - 2.2.04 Glycoprotein
  - 2.2.05 Lipopolysaccharide
  - 2.2.06 Lipoprotein
  - 2.2.07 Phospholipids
  - 2.2.08 Polysaccharides - (cytoplasmic)
  - 2.2.09 Protein modification
  - 2.2.10 Proteins - translation and modification
  - 2.2.11 RNA synthesis, modif., DNA transcrip.
  - 2.2.12 tRNA
- 3.0.0 Metabolism of small molecules
- 3.1.0 Amino acid biosynthesis
  - 3.1.01 Alanine
  - 3.1.02 Arginine
  - 3.1.03 Asparagine
  - 3.1.04 Aspartate
  - 3.1.05 Chorismate
  - 3.1.06 Cysteine
  - 3.1.07 Glutamate
  - 3.1.08 Glutamine
  - 3.1.09 Glycine
  - 3.1.10 Histidine
  - 3.1.11 Isoleucine
  - 3.1.12 Leucine
  - 3.1.13 Lysine
  - 3.1.14 Methionine
  - 3.1.15 Phenylalanine
  - 3.1.16 Proline
  - 3.1.17 Serine
  - 3.1.18 Threonine
  - 3.1.19 Tryptophan
  - 3.1.20 Tyrosine
  - 3.1.21 Valine

**Appendix VIII (cont):**

- 3.2.0 Biosynthesis of cofactors, carriers
  - 3.2.01 Acyl carrier protein (ACP)
  - 3.2.02 Biotin
  - 3.2.03 Cobalamin
  - 3.2.04 Enterochelin
  - 3.2.05 Folic acid
  - 3.2.06 Heme, porphyrin
  - 3.2.07 Lipoate
  - 3.2.08 Menaquinone, ubiquinone
  - 3.2.09 Molybdopterin
  - 3.2.10 Pantothenate
  - 3.2.11 Pyridine nucleotide
  - 3.2.12 Pyridoxine
  - 3.2.13 Riboflavin
  - 3.2.14 Thiamin
  - 3.2.15 Thioredoxin, glutaredoxin, glutathione
  - 3.2.16 biotin carboxyl carrier protein (BCCP)
- 3.3.0 Central intermediary metabolism
  - 3.3.01 2'-Deoxyribonucleotide metabolism
  - 3.3.02 Amino sugars
  - 3.3.03 Entner-Doudoroff
  - 3.3.04 Gluconeogenesis
  - 3.3.05 Glyoxylate bypass
  - 3.3.06 Incorporation metal ions
  - 3.3.07 Misc. glucose metabolism
  - 3.3.08 Misc. glycerol metabolism
  - 3.3.09 Non-oxidative branch, pentose pathway
  - 3.3.10 Nucleotide hydrolysis
  - 3.3.21 other
  - 3.3.11 Nucleotide interconversions
  - 3.3.12 Oligosaccharides
  - 3.3.13 Phosphorus compounds
  - 3.3.14 Polyamine biosynthesis
  - 3.3.15 Pool, multipurpose conversions of intermed. metab.
  - 3.3.16 S-adenosyl methionine
  - 3.3.17 Salvage of nucleosides and nucleotides
  - 3.3.18 Sugar-nucleotide biosynthesis, conversions
  - 3.3.19 Sulfur metabolism
  - 3.3.20 Amino acids
- 3.4.0 Degradation of small molecules
  - 3.4.1 Amines
  - 3.4.2 Amino acids
  - 3.4.3 Carbon compounds
  - 3.4.4 Fatty acids
  - 3.4.5 Other
  - 3.4.0 ATP-proton motive force
- 3.5.0 Energy metabolism, carbon
  - 3.5.1 Aerobic respiration
  - 3.5.2 Anaerobic respiration
  - 3.5.3 Electron transport
  - 3.5.4 Fermentation
  - 3.5.5 Glycolysis
  - 3.5.6 Oxidative branch, pentose pathway
  - 3.5.7 Pyruvate dehydrogenase
  - 3.5.8 TCA cycle
- 3.6.0 Fatty acid biosynthesis
  - 3.6.1 Fatty acid and phosphatidic acid biosynthesis
- 3.7.0 Nucleotide biosynthesis
  - 3.7.1 Purine ribonucleotide biosynthesis
  - 3.7.2 Pyrimidine ribonucleotide biosynthesis
- 4.0.0 Cell envelop
  - 4.1.0 Periplasmic/exported/lipoproteins
  - 4.1.1 Inner membrane
  - 4.1.2 Murein sacculus, peptidoglycan
  - 4.1.3 Outer membrane constituents
  - 4.1.4 Surface polysaccharides & antigens
  - 4.1.5 Surface structures
- 4.2.0 Ribosome constituents
  - 4.2.1 Ribosomal and stable RNAs
  - 4.2.2 Ribosomal proteins - synthesis, modification
  - 4.2.3 Ribosomes - maturation and modification
- 5.0.0 Extrachromosomal
  - 5.1.0 Laterally acquired elements
    - 5.1.1 Colicin-related functions
    - 5.1.2 Phage-related functions and prophages
    - 5.1.5 Pathogenicity island-related function
    - 5.1.3 Plasmid-related functions
    - 5.1.4 Transposon-related functions
- 6.0.0 Global functions
  - 6.1.1 Global regulatory functions
- 7.0.0 Not classified (included putative assignments)



## Appendix IX: List of colour codes

- 0** (white) - Pathogenicity/Adaptation/Chaperones
- 1** (dark grey) - energy metabolism (glycolysis, electron transport etc.)
- 2** (red) - Information transfer (transcription/translation + DNA/RNA modification)
- 3** (dark green) - Surface (IM, OM, secreted, surface structures)
- 4** (dark blue) - Stable RNA
- 5** (Sky blue) - Degradation of large molecules
- 6** (dark pink) - Degradation of small molecules
- 7** (yellow) - Central/intermediary/miscellaneous metabolism
- 8** (light green) - Unknown
- 9** (light blue) - Regulators
- 10** (orange) - Conserved hypo
- 11** (brown) - Pseudogenes and partial genes (remnants)
- 12** (light pink) - Phage/IS elements
- 13** (light grey) - Some misc. information e.g. Prosite, but no function

## Appendix X: List of degenerate nucleotide value/IUB Base Codes.

**R = A or G**

**S = G or C**

**B = C, G or T**

**Y = C or T**

**W = A or T**

**D = A, G or T**

**K = G or T**

**N = A, C, G or T**

**H = A, C or T**

**M = A or C**

**V = A, C or G**

## Appendix XI: Splice site information

Gene	No.	Exon	Intron	Exon	Size (bp)
41-3	1	GAA	<b>GTACACA</b> . . CCTTCTTTTCCATATTTAG	CAA	152
	2	AAT	<b>GTTAAAA</b> . . . TTTT TTTT TTTT TAACTTAG	CCG	208
	3	GAG	<b>GTAAGAA</b> . . . ATTCATTATATATTTATAG	GGA	86
	4	TCG	<b>GTATGGA</b> . . . TTTTGAAATACTTCCCTCAG	TTA	152
	5	ACT	<b>GTAATAT</b> . . TTTT TTTT TTTT TATTCCTAG	ATG	112
	6	CAG	<b>GTAATA</b> . . ATAATGACATTTTGATACAG	ATT	120
	7	AAT	<b>GTACATT</b> . . TTATTTTATTTATTTATAG	AAA	81
	8	TAG	<b>GTATTTG</b> . . ATATTTTACTTATGATAG	TTA	96
RhopH3	1	AGG	<b>GTAATAT</b> . . TTTATTTTATTTT TTTTA	TTT	150
	2	GGA	<b>GTAAGAG</b> . . TTTTATTATTTTATTGTAG	TCC	442
	3	GGA	<b>GTAAGAG</b> . . TTTTATTATTTTATTGTAG	TCC	199
	4	CAG	<b>GTAYGCT</b> . . TTTAATTTT TTTTCTTCA	TCA	160
	5	AAA	<b>GTAAGAA</b> . . TATTTT TTTTACAATTTTAG	TTC	206
	6	AAG	<b>GTAAGAAG</b> . . TTTT TTTT TTTT TGTTCAG	TTT	142
RNA pol III	1	CAG	<b>GTACATA</b> . . TTTT TTTT TTTT TTTTAG	GTG	158
	2	CAA	<b>GTAATTA</b> . . TATATTTTATTTT TCTTAG	GTT	113
	3	TAC	<b>GTTAGTT</b> . . TTTT TTTT TTTT TTTTAG	TGG	169
	4	ATT	<b>GTAAGTT</b> . . TATTTT TTTT TTTT TTTTAG	TGA	112
SERA	1	TGT	<b>GTAAGAA</b> . . TTGTCATTATTTT TTTTAG	GTG	158
	2	AAA	<b>GTATAAA</b> . . TTTATTTATTTT TTTTAG	ATA	175
	3	CAG	<b>GTAATA</b> . . TTTTAATTTT TTTGTTTAG	AAA	129
SERP H	1	CTG	<b>GTTTGTC</b> . . CATATATTTCTTTATTTTAG	ATA	345
	2	AGA	<b>GTAATA</b> . . TTTCTTATATTTTCTTTTAG	GTG	92
	3	CTG	<b>GTTTGTC</b> . . CATATATTTCTTTATTTTAG	ATA	116
Ag15	1	ATG	<b>GTAAGAG</b> . . TATTTT TGATACCTTATAG	AGT	214
	2	AAA	<b>GTAATTA</b> . . CAATCATATTAACACAAAAG	ATG	280
PfGPx	1	GAG	<b>GTATACA</b> . . TTATTATCCCTTGCTTTAG	ATC	208
	2	TCG	<b>GTTAGTA</b> . . TATTTATCATTTT TTTCCAG	ATG	168
Calmodulin	1	GAA	<b>GTAATC</b> . . TTTT TATTTT TCTCATTAG	CTA	480
PfPK1	1	TAG	<b>GTGTGTT</b> . . TCATTACATTTTACCTTAG	GAT	101
MESA	1	TTA	<b>GTAAGTT</b> . . CGTAATATATTTT TTTTAG	GAT	122
Aldolase	1	ATG	<b>GTAAGAA</b> . . TATTTT TATTTT TTTTAG	GCT	452
KAHRP	1	AAC	<b>GTAAGTT</b> . . TTATTT TTTT TCCATATAG	TGC	430
GBPH2	1	TTG	<b>GTATGCC</b> . . TTTGTATTATTTAATTTTAG	AAT	157
GBP	1	TTG	<b>GTATG</b> . . . TGTGTATTGTTATTTTAG	AAT	179
FIRA	1	TGT	<b>GTAAGGA</b> . . TTTT TATATTTT TCTTTAG	CGA	175
GARP	1	AAG	<b>GTAACAA</b> . . TATATGTATTTT TTTTAG	TGC	214

↑  
Donor motif

↑  
Acceptor motif

The splice acceptor and donor sequences for several *P. falciparum* genes: adapted from Coppel and Black(1998). In "Malaria:Parasite Biology, Pathogenesis and Protection", I.W. Sherman (ed.); ASM Press; Washington DC; pp185-202