

# 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

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.

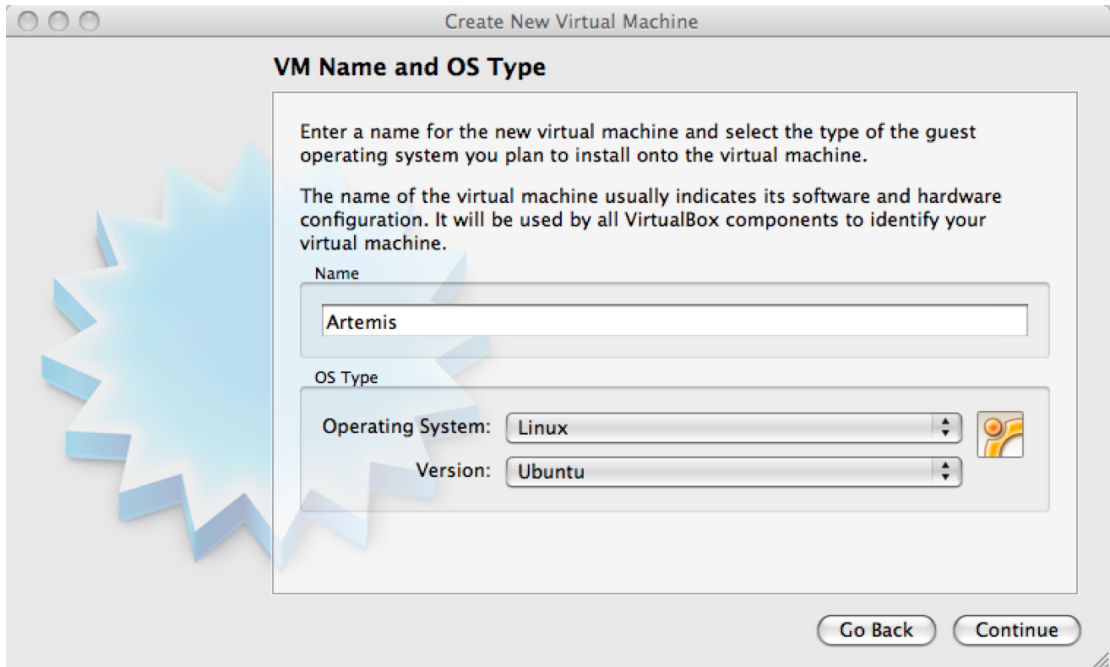
Create a new virtual machine by selecting 'New' from the options at the top. Then fill the boxes in as shown below:

In the first window enter:

Name: **Artemis**

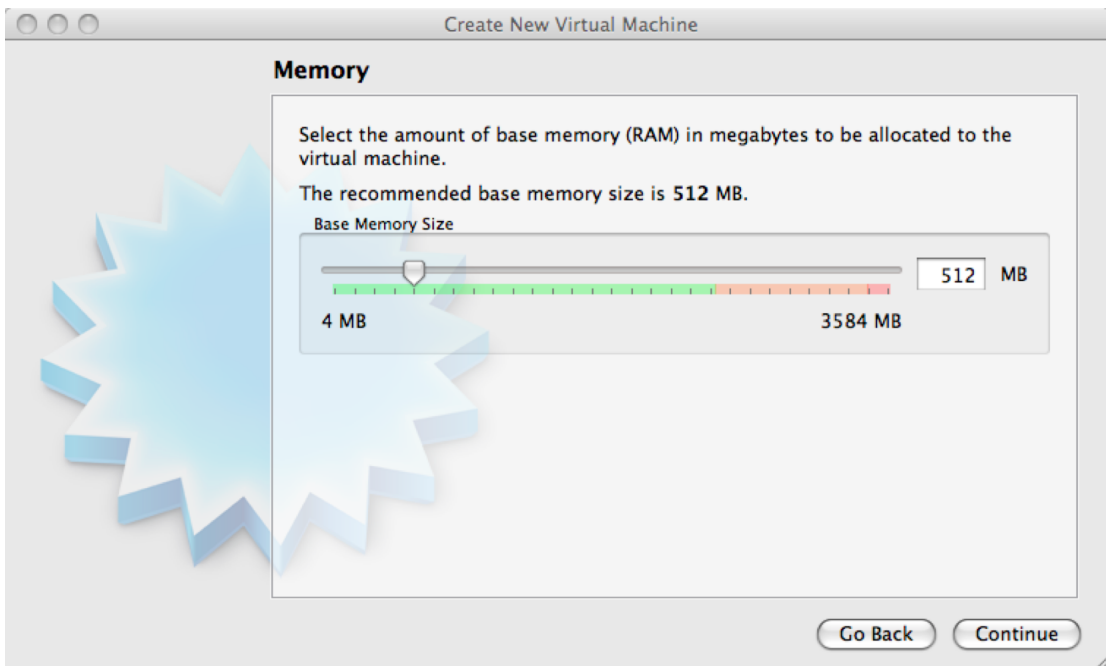
Operating System: **Linux**

Version: **Ubuntu**



Click 'Continue'

In the next window keep the memory default setting (512 MB). You can use more but no more than half the amount of memory on your PC.



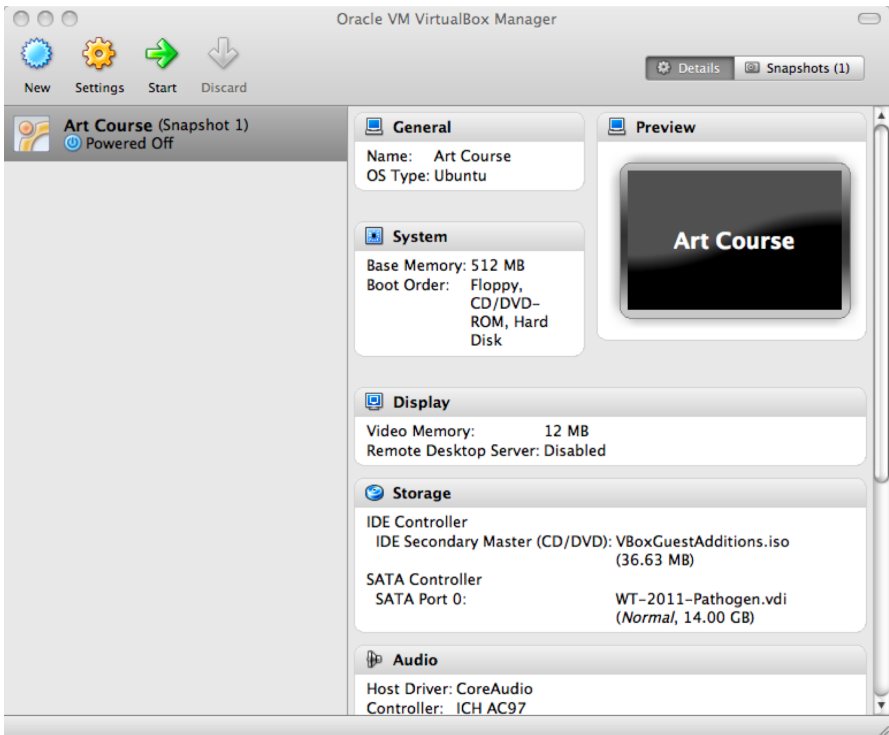
Click 'Continue'.

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 directory. 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.



## **Appendix II: Artemis minimum hardware and software requirements.**

Artemis and ACT will, in general, work well on any standard modern machine and with most common operating systems. It is currently used on many different varieties of UNIX and Linux systems as well as Apple Macintosh and Microsoft Windows systems.

## **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 a 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>)

**Genomes At the EBI**

**Access to Completed Genomes**

The first completed genomes from *viruses*, *phages* and *organelles* were deposited into the EMBL Database in the early 1980's. Since then, molecular biology's shift to obtain the complete sequences of as many genomes as possible combined with major developments in sequencing technology resulted in hundreds of complete genome sequences being added to the database, including *Archaea*, *Bacteria* and *Eukaryota*. These web pages give access to a large number of complete genomes, [help](#) is available to describe the layout.

**Whole Genome Shotgun Sequences (WGS)**

Methods using whole genome shotgun data are used to gain a large amount of genome coverage for an organism. WGS data for a growing number of organisms are being submitted to DDBJ/EMBL/GenBank and are made available via EBI's Sequence Retrieval System (SRS) at <http://srs.ebi.ac.uk> and the EBI FTP server at <ftp://ftp.ebi.ac.uk/pub/databases/embl/wgs/>. [More information about WGS projects...](#)

**Human Draft Genome**

The completion of the human draft genome sequence was announced and published in February 2001 in *Nature* and *Science*. Since the beginning of the Human Genome Project, the international Human Genome Sequencing Consortium has been submitting human draft sequence data to the International Nucleotide Sequence Databases DDBJ/EMBL/GenBank. High-throughput human sequences have been made available to the public immediately via the EMBL Database [high-throughput genome division](#) (HTG), while finished sequences have been included in the [Human division](#) (HUM).

**Genome Annotation and Proteome Analysis**

Click on the Plasmid hyperlink

**Genomes Plasmid**

212 organisms.

Accession numbers of all the entries listed below may be downloaded as a [text file](#) for use in downloading using the [Sequence Version Archive](#).

	Description	Length (bp)	Sequence	
			Plain	HTML
<i>Acetobacter aceti</i>				
1	Acetobacter aceti pAC5	5,123	<a href="#">AF110140</a>	<a href="#">AF110140</a>
<i>Acetobacter pasteurianus</i>				
2	Acetobacter pasteurianus plasmid pAP12875	1,440	<a href="#">U20550</a>	<a href="#">U20550</a>
<i>Acidithiobacillus ferrooxidans</i>				
3a	Acidithiobacillus ferrooxidans pTF4.1 plasmid	4,104	<a href="#">X96982</a>	<a href="#">X96982</a>
3b	Acidithiobacillus ferrooxidans plasmid pTF5	19,792	<a href="#">U73041</a>	<a href="#">U73041</a>
<i>Acinetobacter sp. EB104</i>				
4	Acinetobacter sp. EB104 plasmid pAC450	4,379	<a href="#">AJ311718</a>	<a href="#">AJ311718</a>
<i>Actinobacillus pleuropneumoniae</i>				
5a	Actinobacillus pleuropneumoniae plasmid pMS260	8,124	<a href="#">AB109805</a>	<a href="#">AB109805</a>

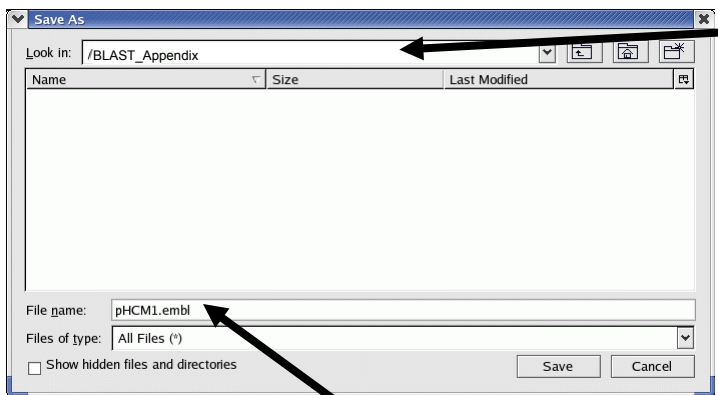
Scroll down the page to the *Salmonella* plasmids

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

Genomes Pages - Plasmid - Microsoft Internet

Address <http://www.ebi.ac.uk/genomes/plasmid.html>

Accession	Plasmid Name	Size (bp)	Accession	Accession	FASTA SRS
<b>Riemerella anatipestifer</b>					
158a	Riemerella anatipestifer plasmid pCFC1	3,966	AF048718	AF048718	4 FASTA SRS
158b	Riemerella anatipestifer plasmid pCFC2	5,609	AF082180	AF082180	3 FASTA SRS
<b>Ruminococcus flavefaciens</b>					
159	Ruminococcus flavefaciens R13e2 cryptic plasmid pBAW301	1,768	U22411	U22411	1 FASTA SRS
<b>Salmonella choleraesuis</b>					
160	Salmonella choleraesuis strain 79500 plasmid pSFD10	4,801	AY048853	AY048853	6 FASTA SRS
<b>Salmonella enterica</b>					
161	Salmonella enterica subsp. enterica serovar Berta plasmid pBERT	4,656	AF025795	AF025795	9 FASTA SRS
162a	Salmonella enterica subsp. enterica serovar Typhi str. CT18 plasmid pHCM1	218,160	AL513383	AL513383	234 FASTA SRS
162b	Salmonella enterica subsp. enterica serovar Typhi str. CT18 plasmid pHCM2	106,516	AL513384	AL513384	132 FASTA SRS
163	Salmonella enterica subsp. enterica serovar Typhimurium plasmid pFPTB1	12,656	AJ634602	AJ634602	6 FASTA SRS
<b>Salmonella enteritidis</b>					
164a	Salmonella enteritidis serovar Enteritidis plasmid pC	5,269	AY079201	AY079201	4 FASTA SRS
164b	Salmonella enteritidis serovar Enteritidis plasmid pK	4,245	AY079200	AY079200	3 FASTA SRS
164c	Salmonella enteritidis serovar Enteritidis plasmid pP	4,301	AY079199	AY079199	3 FASTA SRS
<b>Salmonella typhi</b>					
165a	Salmonella typhi R27 plasmid	180,461	AF250878	AF250878	204 FASTA SRS
165b	Salmonella typhi plasmid R27	38,245	AF105019	AF105019	34 FASTA SRS
<b>Salmonella typhimurium</b>					



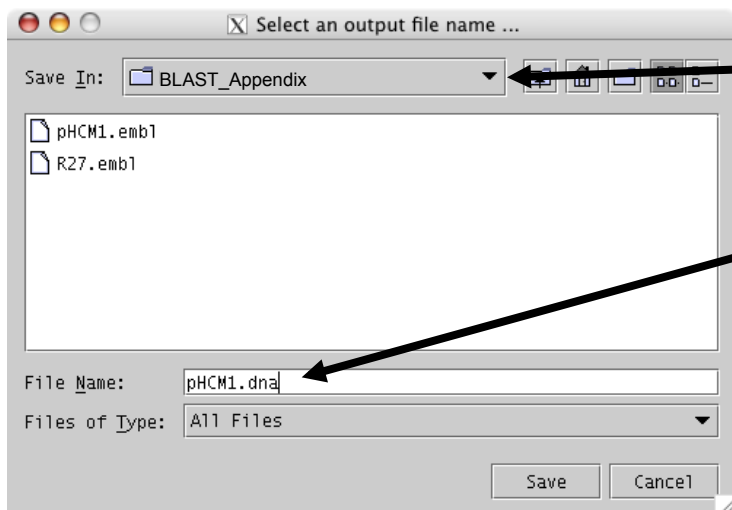
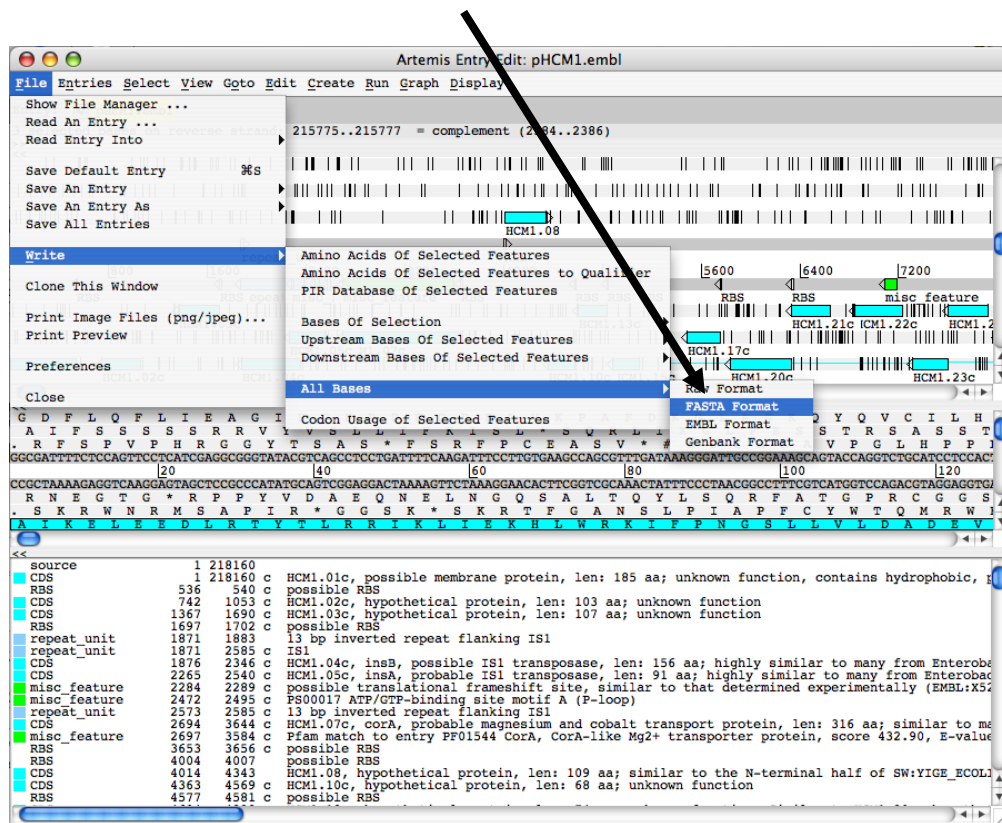
Save the EMBL sequence in a suitable directory. For example: BLAST\_Appendix

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. 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

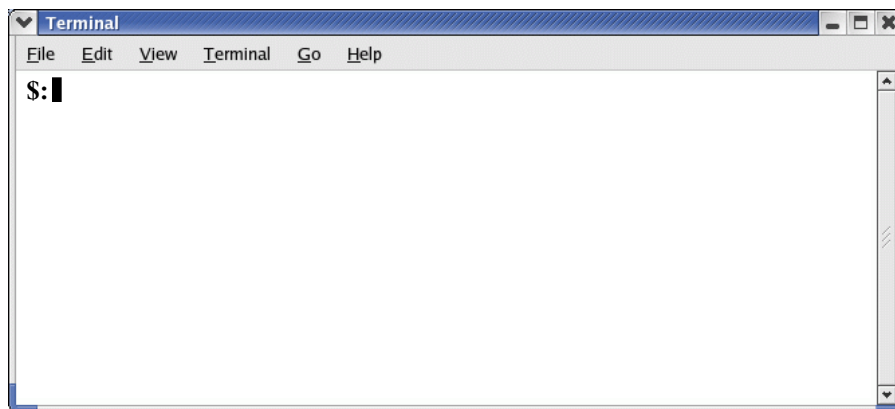
## Running Blast

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.



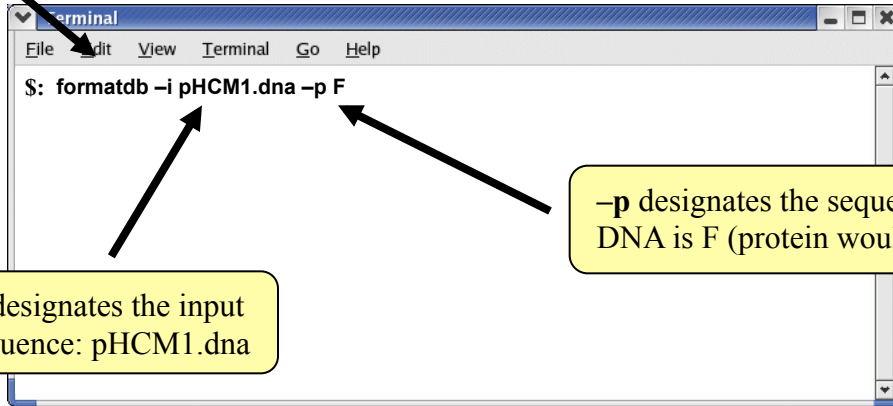
You will treat pHCM1.dna as the **database** sequence and R27.dna as the **query** sequence

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)



```
Terminal
File Edit View Terminal Go Help
$ formatdb -i pHCM1.dna -p F
```

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

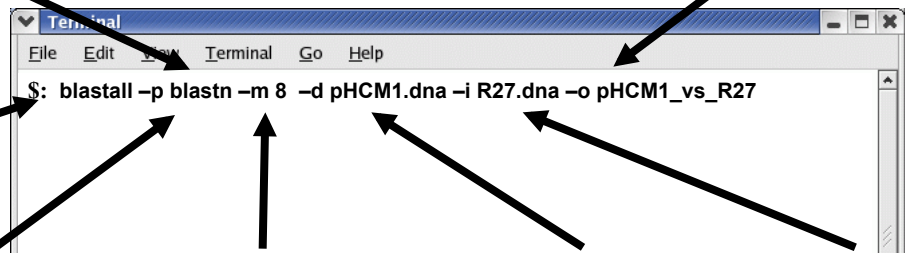
**blastall** is the BLAST program

**-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

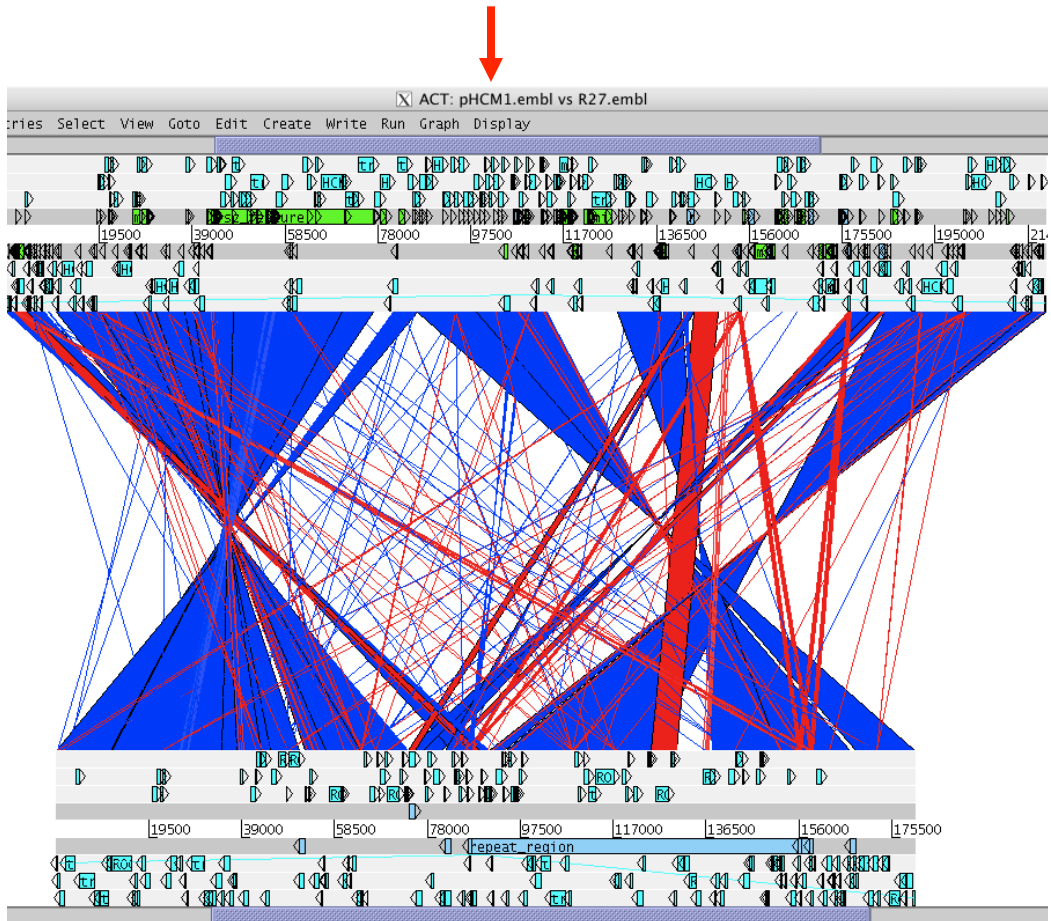
**-i** designates the query sequence: R27.dna



```
Terminal
File Edit View Terminal Go Help
$ blastall -p blastn -m 8 -d pHCM1.dna -i R27.dna -o pHCM1_vs_R27
```



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

EMBL-EBI  
European Bioinformatics Institute

Genomes Pages

188 organisms.

Accession numbers of all the entries listed below may be downloaded as a [text file](#) for use in downloading using the [Sequence Version Archive](#).

	Description	Length (bp)	Sequence		Proteins
			Plain	HTML	
<i>Acinetobacter</i> sp. ADP1 (Description)					
1	<i>Acinetobacter</i> sp. ADP1	3,598,621	CR543861	CR543861	Proteome
<i>Agrobacterium tumefaciens</i>					
2a	<i>Agrobacterium tumefaciens</i> str. C58 (Cereon) chromosome (circular) (254 parts in a CON entry)	2,841,561	AE007869	AE007869	Proteome
2b	<i>Agrobacterium tumefaciens</i> str. C58 (Cereon) chromosome (linear) (187 parts in a CON entry)	2,074,762	AE007870	AE007870	Proteome
2c	<i>Agrobacterium tumefaciens</i> str. C58 (U. Washington) chromosome (circular) (256 parts in a CON entry)	2,841,490	AE008688	AE008688	Proteome
2d	<i>Agrobacterium tumefaciens</i> str. C58 (U. Washington) chromosome (linear) (187 parts in a CON entry)	2,075,560	AE008689	AE008689	Proteome
<i>Anaplasma marginale</i>					
3	<i>Anaplasma marginale</i> str. St. Maries	1,197,687	CP000030	CP000030	Proteome
<i>Aquifex aeolicus</i>					
4	<i>Aquifex aeolicus</i> VF5 (109 parts in a CON entry)	1,551,335	AE000657	AE000657	Proteome
<i>Azoarcus</i> sp. EbN1					

Scroll down the page to the *Staphylococcus aureus* genomes

Genomes Pages - Bacteria - Mozilla

<http://www.ebi.ac.uk/genomes/bacteria.html>

133 *Shewanella oneidensis* MR-1 (457 parts in a CON entry)

134 *Shigella flexneri* 2a str. 301

135 *Shigella flexneri* 2a str. 2457T (16 parts in a CON entry)

136 *Silicibacter pomeroyi* DSS-3

137 *Sinorhizobium meliloti* 1021 (12 parts in a CON entry)

*Staphylococcus aureus*

138 *Staphylococcus aureus* subsp. *aureus* COL

139 *Staphylococcus aureus* subsp. *aureus* MRSA252

140 *Staphylococcus aureus* subsp. *aureus* MSSA476

141 *Staphylococcus aureus* subsp. *aureus* MW2

142 *Staphylococcus aureus* subsp. *aureus* Mu50

143 *Staphylococcus aureus* subsp. *aureus* N315

144 *Staphylococcus epidermidis* ATCC 12228 (9 parts in a CON entry)

145 *Staphylococcus epidermidis* RP62A

*Streptococcus agalactiae*

146 *Streptococcus agalactiae* 2603VR (100 parts in a CON entry)

147 *Streptococcus agalactiae* NEM316 (14 parts in a CON entry)

*Streptococcus mutans*

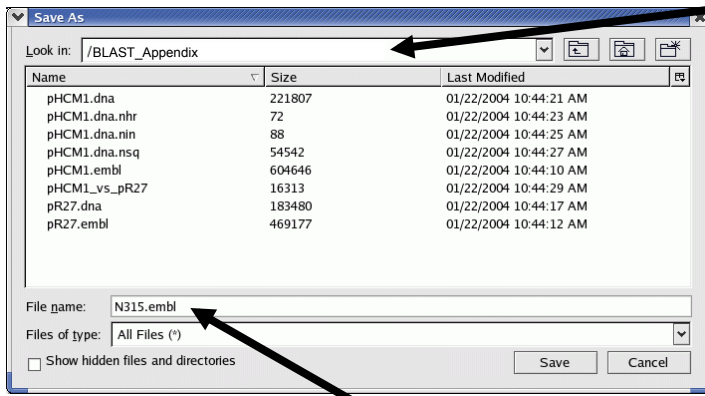
148 *Streptococcus mutans* UA159 (185 parts in a CON entry)

*Streptococcus pneumoniae*

149 *Streptococcus pneumoniae* R6 (184 parts in a CON entry)

150 *Streptococcus pneumoniae* TIGR4 (194 parts in a CON entry)

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

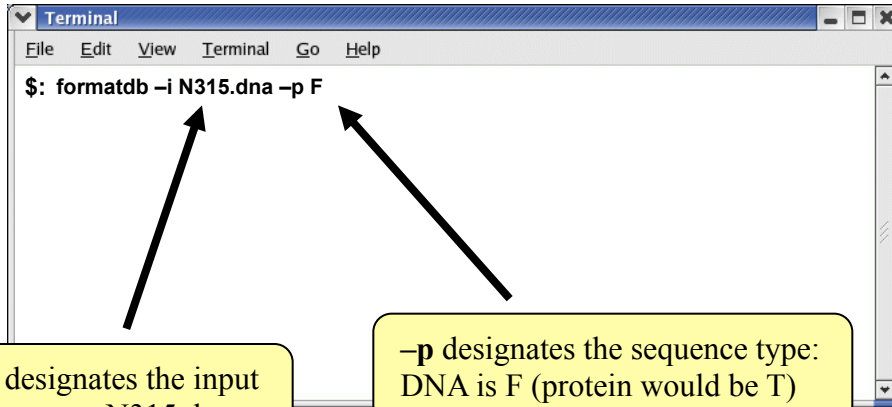
In the previous exercise you used the **blastall** program to run BLASTN on two plasmid sequences. As the genome sequences are larger (~2.8 Mb) you are going to run **megablast**, another program from the NCBI BLAST distribution that can generate comparison files in a format that ACT can read (see Appendix 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**



```
$: formatdb -i N315.dna -p F
```

Diagram description: A terminal window titled 'Terminal' with a menu bar (File, Edit, View, Terminal, Go, Help). The command '\$: formatdb -i N315.dna -p F' is entered. A red arrow points from the instruction box above to the terminal. Two black arrows point from callout boxes to the command: one to '-i' and another to '-p'.

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

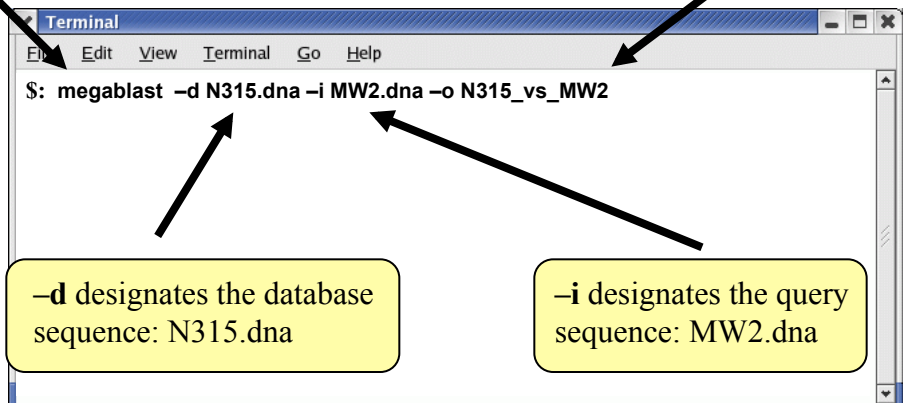
**-p** designates the sequence type: DNA is F (protein would be T)

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



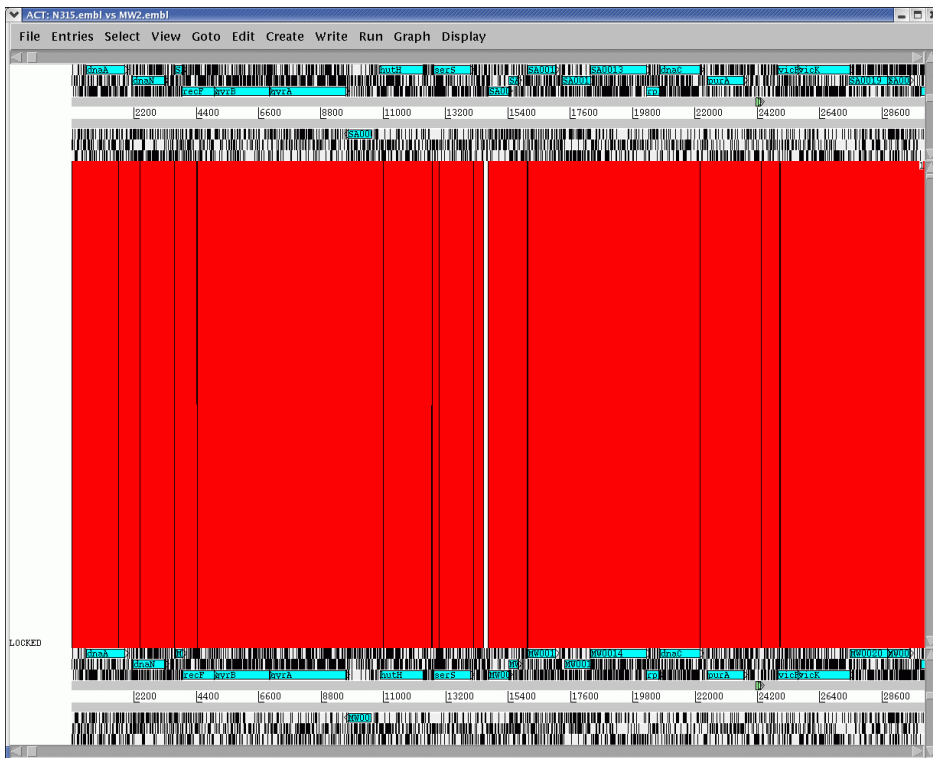
```
$: megablast -d N315.dna -i MW2.dna -o N315_vs_MW2
```

Diagram description: A terminal window titled 'Terminal' with a menu bar (File, Edit, View, Terminal, Go, Help). The command '\$: megablast -d N315.dna -i MW2.dna -o N315\_vs\_MW2' is entered. A red arrow points from the instruction box above to the terminal. Four black arrows point from callout boxes to the command: one to 'megablast', one to '-d', one to '-i', and one to '-o'.

**-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 VI: Useful Web addresses

### Major Public Sequence Repositories

DNA Data Bank of Japan (DDBJ)	<a href="http://www.ddbj.nig.ac.jp">http://www.ddbj.nig.ac.jp</a>
EMBL Nucleotide Sequence Database	<a href="http://www.ebi.ac.uk/embl">http://www.ebi.ac.uk/embl</a>
Genomes at the EBI	<a href="http://www.ebi.ac.uk/genomes">http://www.ebi.ac.uk/genomes</a>
GenBank	<a href="http://www.ncbi.nih.gov/Genbank">http://www.ncbi.nih.gov/Genbank</a>

### Microbial Genome Databases Resources

Sanger Microbial Genomes	<a href="http://www.sanger.ac.uk/Projects/Pathogens">http://www.sanger.ac.uk/Projects/Pathogens</a>
GeneDB	<a href="http://www.genedb.org">http://www.genedb.org</a>
Institute Pasteur GenoList databases <i>Including: SubtiList, Colbri, TubercuList, Leproma, PyloriGene, MypuList, ListiList, CandidaDB.</i>	<a href="http://genolist.pasteur.fr">http://genolist.pasteur.fr</a>
Pseudomonas Genome Database	<a href="http://www.pseudomonas.com">http://www.pseudomonas.com</a>
Clusters of Orthologous Groups of proteins (COGs)	<a href="http://www.ncbi.nlm.nih.gov/COG">http://www.ncbi.nlm.nih.gov/COG</a>
ScoDB ( <i>S. coelicolor</i> database)	<a href="http://streptomyces.org.uk">http://streptomyces.org.uk</a>
GenProtEC	<a href="http://genprotec.mbl.edu">http://genprotec.mbl.edu</a>

### Protein Motif Databases

Prosite	<a href="http://www.expasy.ch/prosite/">http://www.expasy.ch/prosite/</a>
Pfam	<a href="http://pfam.sanger.ac.uk">http://pfam.sanger.ac.uk</a>
BLOCKS	<a href="http://blocks.fhcrc.org">http://blocks.fhcrc.org</a>
InterPro	<a href="http://www.ebi.ac.uk/interpro/">http://www.ebi.ac.uk/interpro/</a>
PRINTS	<a href="http://umber.sbs.man.ac.uk/dbbrowser/PRINTS/">http://umber.sbs.man.ac.uk/dbbrowser/PRINTS/</a>
SMART	<a href="http://smart.embl-heidelberg.de">http://smart.embl-heidelberg.de</a>

### Protein feature prediction tools

TMHMM Transmembrane helices prediction	<a href="http://www.cbs.dtu.dk/services/TMHMM-2.0/">http://www.cbs.dtu.dk/services/TMHMM-2.0/</a>
SignalP Prediction Server	<a href="http://www.cbs.dtu.dk/services/SignalP/">http://www.cbs.dtu.dk/services/SignalP/</a>
PSORT protein prediction	<a href="http://psort.ims.u-tokyo.ac.jp/form.html">http://psort.ims.u-tokyo.ac.jp/form.html</a>

### Metabolic Pathways and Cellular Regulation

EcoCyc	<a href="http://ecocyc.org/">http://ecocyc.org/</a>
ENZYME	<a href="http://www.expasy.ch/enzyme/">http://www.expasy.ch/enzyme/</a>
Kyoto Encyclopedia of Genes and Genomes (KEGG)	<a href="http://www.genome.ad.jp/kegg">http://www.genome.ad.jp/kegg</a>
MetaCyc	<a href="http://metacyc.org/">http://metacyc.org/</a>

### Miscellaneous sites

NCBI BLAST website	<a href="http://www.ncbi.nlm.nih.gov/BLAST/">http://www.ncbi.nlm.nih.gov/BLAST/</a>
EBI FASTA website	<a href="http://www.ebi.ac.uk/fasta33/index.html">http://www.ebi.ac.uk/fasta33/index.html</a>
The tmRNA website	<a href="http://www.indiana.edu/~tmrna/">http://www.indiana.edu/~tmrna/</a>
tRNAscan-SE Search Server	<a href="http://selab.janelia.org/tRNAscan-SE/">http://selab.janelia.org/tRNAscan-SE/</a>
Rfam	<a href="http://rfam.sanger.ac.uk/">http://rfam.sanger.ac.uk/</a>
Codon usage database	<a href="http://www.kazusa.or.jp/codon/">http://www.kazusa.or.jp/codon/</a>
GO Gene Ontology Consortium	<a href="http://www.geneontology.org/">http://www.geneontology.org/</a>
Artemis homepage	<a href="http://www.sanger.ac.uk/Software/Artemis/">http://www.sanger.ac.uk/Software/Artemis/</a>
ACT homepage	<a href="http://www.sanger.ac.uk/Software/ACT/">http://www.sanger.ac.uk/Software/ACT/</a>
WebACT	<a href="http://www.webact.org/WebACT/home">http://www.webact.org/WebACT/home</a>
Double ACT	<a href="http://www.hpa-bioinfotools.org.uk/pise/double_act.html">http://www.hpa-bioinfotools.org.uk/pise/double_act.html</a>
Glimmer	<a href="http://cbbcb.umd.edu/software/glimmer/">http://cbbcb.umd.edu/software/glimmer/</a>
EasyGene	<a href="http://www.cbs.dtu.dk/services/EasyGene/">http://www.cbs.dtu.dk/services/EasyGene/</a>
String	<a href="http://string.embl.de">http://string.embl.de</a>
EMBOSS	<a href="http://emboss.sourceforge.net/">http://emboss.sourceforge.net/</a>

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

This scheme was adapted for in-house use from the Monica Riley's 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

2.1.3 Degradation of polysaccharides

2.1.2 Degradation of RNA

2.1.4 Degradation of proteins, peptides, glycoproteins

2.2.0 Macromolecule synthesis, modification

2.2.01 Amino acyl tRNA synthesis; tRNA modification

2.2.07 Phospholipids

2.2.02 Basic proteins - synthesis, modification

2.2.08 Polysaccharides - (cytoplasmic)

2.2.03 DNA - replication, repair, restriction./modification

2.2.09 Protein modification

2.2.04 Glycoprotein

2.2.10 Proteins - translation and modification

2.2.05 Lipopolysaccharide

2.2.11 RNA synthesis, modif., DNA transcrip.

2.2.06 Lipoprotein

2.2.12 tRNA

3.0.0 Metabolism of small molecules

3.1.0 Amino acid biosynthesis

3.1.01 Alanine

3.1.08 Glutamine

3.1.15 Phenylalanine

3.1.02 Arginine

3.1.09 Glycine

3.1.16 Proline

3.1.03 Asparagine

3.1.10 Histidine

3.1.17 Serine

3.1.04 Aspartate

3.1.11 Isoleucine 3.1.18 Threonine

3.1.05 Chorismate

3.1.12 Leucine

3.1.19 Tryptophan

3.1.06 Cysteine

3.1.13 Lysine

3.1.20 Tyrosine

3.1.07 Glutamate

3.1.14 Methionine

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
  - 3.3.11 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
  - 3.3.21 other
- 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
    - 5.1.3 Plasmid-related functions
    - 5.1.4 Transposon-related functions
    - 5.1.5 Pathogenicity island-related function
- 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 Xa: 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**

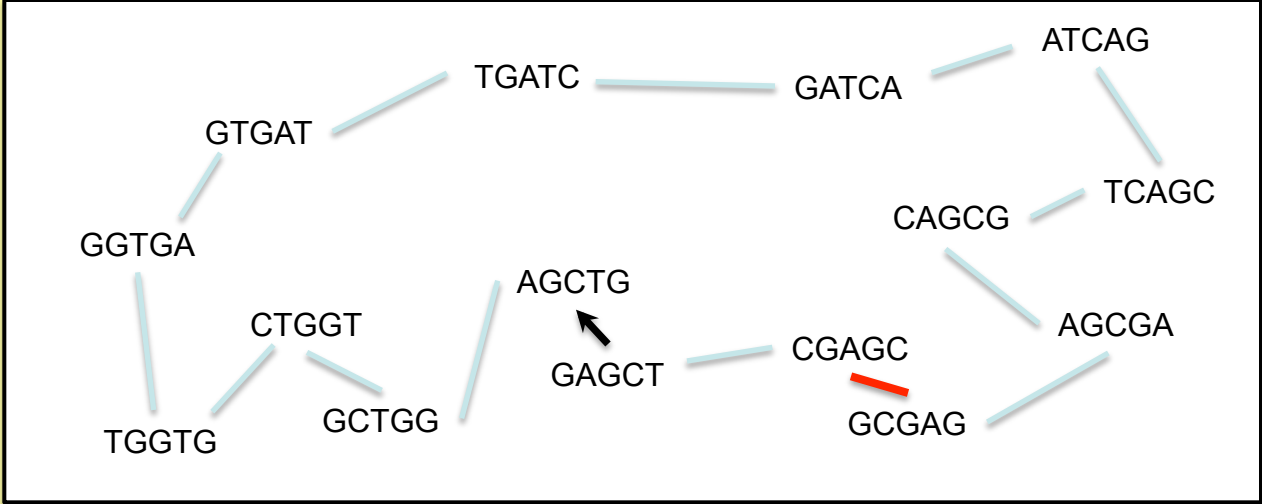
**M = A or C**

**V = A, C or G**

# Appendix X - Assembly

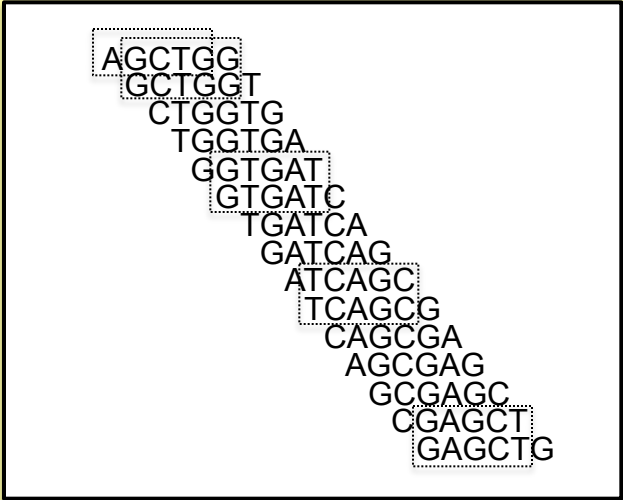
Here we present the solution for the de Bruijn graph exercise, as well as the code for the PERL scripts we mentioned in the assembly module.

The exercise is on page 2 of the assembly module. The first step of the solution would be to generate all the k-mers from the reads. For example the k-mers of length 5 for the read GCGAGC are GCGAG and CGAGC. Those two k-mers will be nodes in the de Bruijn graph, and moreover, will be connected (bold red edge). Doing this for all the reads, generates the following graph:



It is not always easy and might be confusing which node to connect. Remember, the concept of k-mers and the de Bruijn graph are needed to be able to process the large amount of short reads generated by the sequencing machines.

The graph can also be represented as a multiple alignment, as shown on the right hand site. All reads are aligned against each other. The dotted boxes are examples of k-mers.



Now just follow the path through the graph. Starting at the arrow, the first k-mer is GAGCT, so this would be the start of our contig. The graph indicates the next k-mer AGCTG. So we add a G to the contig. The next k-mer is GCTGG. The new letter is another G. Doing this for the whole graph, we get: **GAGCTGGTGATCAGC**. As you see, the graph is circular. So depending where you start, you get a different contig! If you do a six frame translation, you might see which is a good starting point for the contig.

## PERL: Find read pairs that map too far apart

For some applications it would be useful to know whether read pairs map too far apart or whether they don't map pointing to each other. This could be an indication of mis-assemblies, but also duplications or rearrangements, which are looking for when comparing sequences of different strains.

To find read pairs (RPs) that map too far apart we just need columns 2 and 9 from the BAM file (mapping flag and insert size), and a PERL one-liner. We successively make the query more and more complex, until we find the mis-assembly. Please keep in mind that this is advanced programming! It should give you an idea how useful programming could be.

Assuming your BAM file is called IT\_onDenovo.bam and you want to list RPs that map more than 2000bp apart:

```
$ samtools view IT_onDenovo.bam | perl -nle 'my ($read,$flag,$ref,$pos,$mappingQual,$cigar,$mateRef,$matePos,$insertSize,$seq,$seqQual,$other)=split(/\t/); if($insertSize>2000){print}' | head
```

Here is also a shorter version, using an array (not as readable):

```
$ samtools view IT_onDenovo.bam | perl -nle 'my @ar=split(/\t/); if($ar[8]>2000){print}' | head
```

There is a lot of output. Many read pairs map all over the place. We would like to bin those into chunks of 1kb, and then list of the most abundant:

```
$ samtools view IT_onDenovo.bam | perl -nle 'my ($read,$flag,$ref,$pos,$mappingQual,$cigar,$mateRef,$matePos,$insertSize,$seq,$seqQual,$other)=split(/\t/); if($insertSize>2000){print int($pos/1000)."\t".int($matePos/1000)}' | sort | uniq -c | sort -rn | head
```

This does look more complex! In the output the first column is the number of RPs that connect the first bin (2<sup>nd</sup> column) with the second bin (3<sup>rd</sup> column). For example 418 1360 1373 means that 481 RPs connect the region 1360000-1361000 of genome with the regions 1373000-13731000 of the genome. The list shows us that in the subtelomeric regions many RP map far apart!

The following command ignores the subtelomeric ends, by excluding 75kb at each end.

```
$ samtools view IT_onDenovo.bam | perl -nle 'my ($read,$flag,$ref,$pos,$mappingQual,$cigar,$mateRef,$matePos,$insertSize,$seq,$seqQual,$other)=split(/\t/); if($insertSize>10000 && $pos>75000 && $matePos < 1300000){print int($pos/1000)."\t".int($matePos/1000)}' | sort | uniq -c | sort -nr
```

The third line 21 81 323 shows us our mis-assembly. What are the other entries?

We are fully aware that this is a quite complex piece of code, and just used as a one liner. It uses the LINUX commands sort and uniq. But keep in mind that this command can find you all mate pairs mapping too far apart for any bam file (if you adjust the insertSize parameter for your data)!

# Getting non mapping reads and their mates

Here is an example of how to get the mates of non mapping reads. It is a good example of PERL one-liners.

First we are going to get reads that don't map with PERL. The original command is:

```
$ samtools view -f 0X4 IT.Chr5.bam | head
```

In PERL this would be:

```
$ samtools view IT.Chr5.bam | perl -nle 'my ($read, $flag)=split(/\t/); if ($flag & 0x4) {print }' | head
```

Now we need to get the reads where the mate is not mapped. Looking at the samtools manual:

Bit	Description
0x1	template having multiple segments in sequencing
0x2	each segment properly aligned according to the aligner
0x4	segment unmapped
0x8	next segment in the template unmapped
0x10	SEQ being reverse complemented
0x20	SEQ of the next segment in the template being reversed
0x40	the first segment in the template
0x80	the last segment in the template
0x100	secondary alignment
0x200	not passing quality controls
0x400	PCR or optical duplicate

The 0x8 tells if the mate pair is not mapped. So if the read is not mapping (0x4) or the mate is not mapping (0x8) then print the sam line into a file:

```
$ samtools view IT.Chr5.bam | perl -nle 'my ($read, $flag)=split(/\t/); if ($flag & 0x4 or $flag & 0x8) {print }' | sort > NonmappingReadsPlusmate.sam
```

This file can now be used in VELVET for *de novo* assembly as explained in the Assembly module.

We hope that this illustrates the power of PERL one-liners!

## 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 TTTTAAACTTAG	CCG	208
	3	GAG	<b>GTAAGAA</b> . . . ATTCATTATATATTTATAG	GGA	86
	4	TCG	<b>GTATGGA</b> . . . TTTTGAAATACTTCCTCAG	TTA	152
	5	ACT	<b>GTAATAT</b> . . TTTT TTTT TTTTATTTCCCTAG	ATG	112
	6	CAG	<b>GTAATA</b> . . . ATAATGACATTTTGATACAG	ATT	120
	7	AAT	<b>GTACATT</b> . . TTATTTT TTTTATTTATTTATAG	AAA	81
	8	TAG	<b>GTATTTG</b> . . ATATTTT TTTTACTTATGATAG	TTA	96
RhopH3	1	AGG	<b>GTAATAT</b> . . TTTATTTTATTTTTTTTTTA	TTT	150
	2	GGA	<b>GTAAGAG</b> . . TTTT TATTATTTTATTTGTTAG	TCC	442
	3	GGA	<b>GTAAGAG</b> . . TTTT TATTATTTTATTTGTTAG	TCC	199
	4	CAG	<b>GTAYGCT</b> . . TTTAATTTT TTTTTCCTTCA	TCA	160
	5	AAA	<b>GTAAGAA</b> . . TATTTT TTTTACAATTTTAG	TTC	206
	6	AAG	<b>GTAAGA</b> . . TTTT TTTT TTTTGTTCAG	TTT	142
RNA pol III	1	CAG	<b>GTACATA</b> . . TTTT TTTT TTTT TTTT TTTAG	GTG	158
	2	CAA	<b>GTAATTA</b> . . TATATTTTATTTTCTTAG	GTT	113
	3	TAC	<b>GTTAGTT</b> . . TTTT TTTT TTTT TTTT TTTAG	TGG	169
	4	ATT	<b>GTAAGTT</b> . . TATTTT TTTT TTTT TTTT TTTAG	TGA	112
SERA	1	TGT	<b>GTAAGAA</b> . . TTGTCATTATTTTTTTTAG	GTG	158
	2	AAA	<b>GTATAAA</b> . . TTTATTTATTTTTTTTAG	ATA	175
	3	CAG	<b>GTAATA</b> . . TTTTAATTTT TTTTGTTTAG	AAA	129
SERP H	1	CTG	<b>GTTTGTC</b> . . CATATATTTCTTTATTTTAG	ATA	345
	2	AGA	<b>GTAAAAA</b> . . TTTCTTATATTTCTTTTAG	GTG	92
	3	CTG	<b>GTTTGTC</b> . . CATATATTTCTTTATTTTAG	ATA	116
Ag15	1	ATG	<b>GTAAGAG</b> . . TATTTT TGATACCTTTATAG	AGT	214
	2	AAA	<b>GTAATTA</b> . . CAATCATATTAACACAAAAG	ATG	280
PfGPx	1	GAG	<b>GTATACA</b> . . TTATTTATTCCTTGCTTTAG	ATC	208
	2	TCG	<b>GTTAGTA</b> . . TATTTATCATTTTTCCTAG	ATG	168
Calmodulin	1	GAA	<b>GTAATC</b> . . TTTT TTTT TTTTCTCATTAG	CTA	480
PfPK1	1	TAG	<b>GTGTGTT</b> . . TCATTACATTTTACCTTAG	GAT	101
MESA	1	TTA	<b>GTAAGTT</b> . . CGTAATATATTTTTTTTAG	GAT	122
Aldolase	1	ATG	<b>GTAAGAA</b> . . TATTTT TATATTTT TTTTAG	GCT	452
KAHRP	1	AAC	<b>GTAAGTT</b> . . TTATTT TTTT TTTTCATATAG	TGC	430
GBPH2	1	TTG	<b>GTATGCC</b> . . TTTGTATTATTTAATTTTAG	AAT	157
GBP	1	TTG	<b>GTATG</b> . . . . TGTGTATTGTTTATTTTAG	AAT	179
FIRA	1	TGT	<b>GTAAGGA</b> . . TTTT TATATTTT TTTCTTAG	CGA	175
GARP	1	AAG	<b>GTAACAA</b> . . TATATGTATTTTTTTTAG	TGC	214

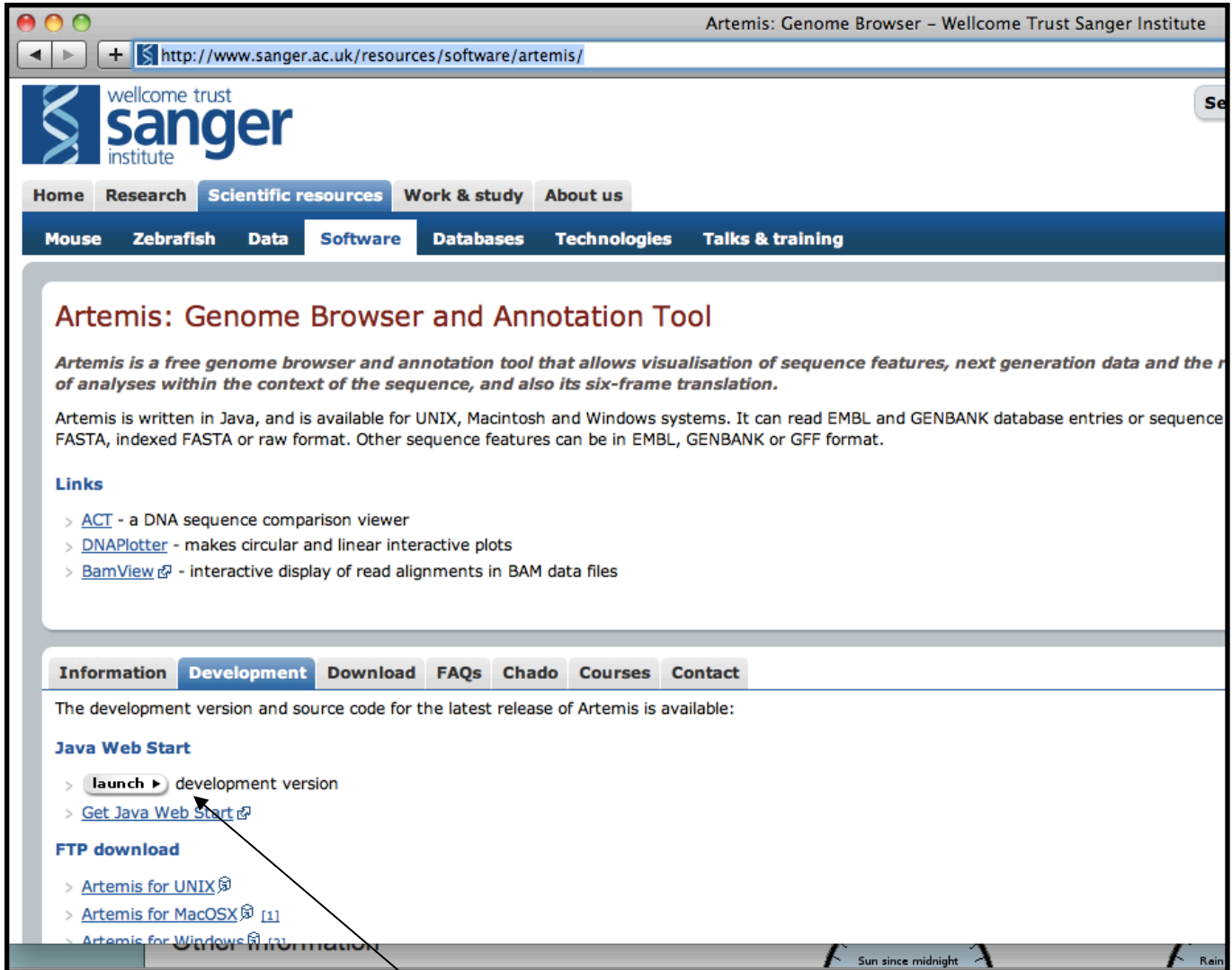


The splice acceptor and donor sequences for several *P. falciparum* genes: adapted from Coppel and Black(1998). In "Malaria:Parasite Biology, Pathogenesis and Protection", I.W. Sherman (ed.); ASM Press; Washington DC; pp185-202

## Appendix XII Running Artemis from the Web

To work this Artemis you don't necessary have to work with it from the VM. It can be run from the web:

<http://www.sanger.ac.uk/resources/software/artemis/>



Click on 'launch', accept and wait a bit...  
Next you can load sequence, bam files etc.

## Appendix XIII Exploring the Sequence Read Archive

### Exercise

You can access the sequence read archive through the following sites:

<http://www.ncbi.nlm.nih.gov/sra>

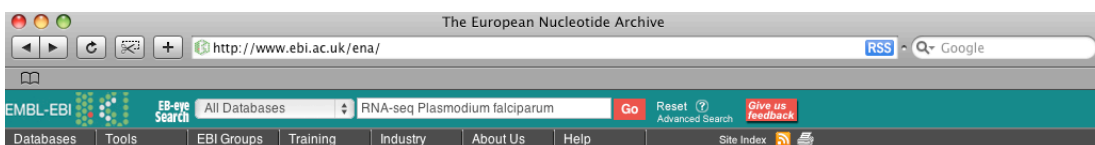
<http://www.ebi.ac.uk/ena>

It's possible to download transcriptome data and other next generation sequencing data as follows. These instructions are given in the form of an exercise to help make it more interesting:

Go to the following website: <http://www.ebi.ac.uk/ena> and type in the search box:

RNA-seq, Plasmodium falciparum

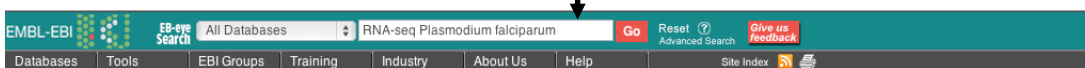
Now follow the step-by-step instructions to download this data.



EBI > Databases > Nucleotide > The European Nucleotide Archive

#### The European Nucleotide Archive

Documentation coming soon, please refer to [documentation](#) relating to assembled sequence and annotation, [information](#) on the Sequence Read Archive (SRA) and the [European Nucleotide Archive Team](#) web pages.

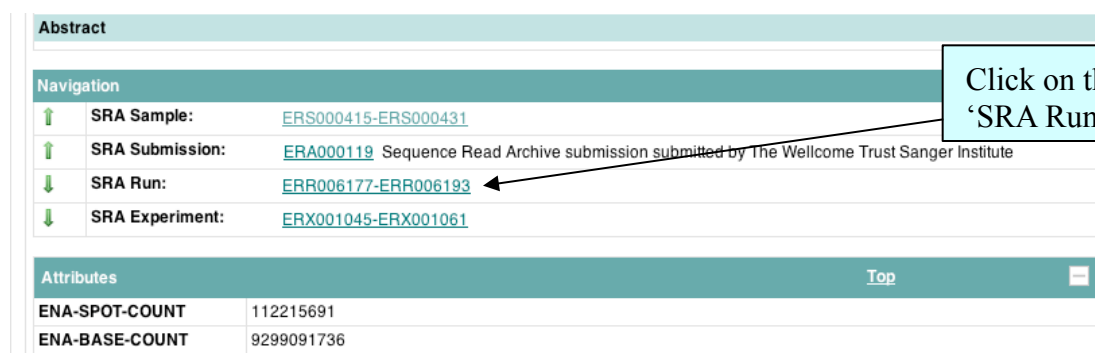


Search for *RNA-seq Plasmodium falciparum* in *All the EBI*

[Expand all](#) [Collapse all](#)

<a href="#">Genomes</a>	0	<a href="#">Molecular Interactions</a>	0
<a href="#">Nucleotide Sequences</a>	1	<a href="#">Reactions &amp; Pathways</a>	0
<a href="#">Protein Sequences</a>	0	<a href="#">Protein Families</a>	0
<a href="#">Macromolecular Structures</a>	0	<a href="#">Enzymes</a>	0
<a href="#">Small molecules</a>	0	<a href="#">Literature</a>	0
<a href="#">Gene Expression</a>	0	<a href="#">Ontologies</a>	0
		<a href="#">EBI Web Site</a>	0

Click on 'Nucleotide Sequences'



Click on the link 'SRA Run'



SRA Run: ERR006185 : 

SRA Run: ERR006186 : 

A list of all the RNA experiments will come up. Click on the red arrow to expand the window.



EMBL-EBI

EB-eye Search

All Databases

Enter Text Here

Go

Reset

Advanced Search

Give us feedback

Databases

Tools

EBI Groups

Training

Industry

About Us

Help

Site Index

Documentation

People

Contact

SRA Experiment Record : ERX001045 : Illumina Genome Analyzer II sequencing of Plasmodium falciparum 3D7 16 hr cDNA-50

View: [XML](#)

Download: [XML](#)

Study : ERP000069 : Plasmodium falciparum RNA-Seq in Blood Stage

More details

Sample : ERS000416 :

More details

Taxonomic classification

Organism

Plasmodium falciparum 3D7

Taxonomic identifier

36329

Library

Name :

16\_hr\_3D7\_cDNA-50\_1

Source :

NON GENOMIC

Selection :

cDNA

Orientation :

Forward

Paired :

Yes

Nominal length :

150

Nominal sdev :

0.0

Spot descriptor

Class :

Application Read

Read 0: Type :

Forward

Base coordinates :

1

Class :

Application Read

Read 1: Type :

Reverse

Base coordinates :

55

Platform

ILLUMINA : Illumina Genome Analyzer II

Processing

Base caller :

Solexa primary analysis

Quality scorer :

Solexa primary analysis

Runs

ERR006186

Submission

ERA000119

Click on 'Runs'



EMBL-EBI

EB-eye Search

All Databases

Enter Text Here

Go

Reset

Advanced Search

Give us feedback

Databases

Tools

EBI Groups

Training

Industry

About Us

Help

Site Index

Documentation

People

Contact

SRA Run Record : ERR006186 : Sequence Read Archive run by The Wellcome Trust Sanger Institute

View: [XML](#)

Download: [XML](#)

Study : ERP000069 : Plasmodium falciparum RNA-Seq in Blood Stage

More details

Sample : ERS000416 :

More details

Taxonomic classification

Organism

Plasmodium falciparum 3D7

Taxonomic identifier

36329

Submitter

The Wellcome Trust Sanger Institute

Experiment

ERX001045

Submission

ERA000119

\* For Aspera download, please [download and install Aspera Connect](#)

Submitted files

ftp

Download

Generated files

ftp

ERR006186\_1.fastq.gz

ERR006186\_2.fastq.gz

Now you can download the RNA-Seq data.



## Appendix XIV

Here is compilations of the programs that we find useful.

### Artemis & Act

<http://www.sanger.ac.uk/resources/software/artemis/#development>

<http://www.sanger.ac.uk/resources/software/act/#development>

### Mapper

SMALT: [http://www.sanger.ac.uk/resources/software/smalt/#t\\_2](http://www.sanger.ac.uk/resources/software/smalt/#t_2)

BWA: <http://sourceforge.net/projects/bio-bwa/files/>

MAQ: <http://sourceforge.net/projects/maq/files/maq/0.6.6/>

BOWTIE: <http://sourceforge.net/projects/bowtie-bio/files/bowtie/0.12.7/>

Top Hat: <http://tophat.cbcb.umd.edu/downloads/>

### SAMTOOLS & BCFtools

svn co <https://samtools.svn.sourceforge.net/svnroot/samtools/trunk/samtools>

### Assemblers

Velvet: <http://www.ebi.ac.uk/~zerbino/velvet/>

ABYSS: <http://www.bcgsc.ca/platform/bioinfo/software/abyss>

SOAPdenovo: <http://soap.genomics.org.cn/soapdenovo.html>

### Tools for automatic finishing / Annotation transfer

ABACAS: <http://sourceforge.net/projects/abacas/files/>

IMAGE: <http://sourceforge.net/projects/image2/>

iCORN: <http://sourceforge.net/projects/icorn/files/>

RATT: svn co <https://ratt.svn.sourceforge.net/svnroot/ratt> ratt

PAGIT: <http://www.sanger.ac.uk/resources/software/pagit/>