

Module 2

Artemis

Introduction

Artemis is a DNA viewer and annotation tool, free to download and use from the Sanger Institute (Rutherford *et al.*, 2000). The program allows the user to view a range of files, from simple sequence files (e.g. fasta format) to EMBL/Genbank entries, as well as the results of sequence analyses, in a highly interactive and intuitive graphical format. Artemis is routinely used for annotation and analysis of both prokaryotic and eukaryotic genomes, and can also be used to visualize mapped data from next generation sequencing. Several types of information can be viewed simultaneously within different contexts. For example, Artemis gives you two views of the same genome region, so you can zoom in to inspect detailed DNA sequence motifs, and also zoom out to view local gene architecture (e.g. operons), or even an entire chromosome or genome, all within one screen. It is also possible to perform analyses within Artemis and save the output for future reference.

Aims

The aim of this Module is for you to become familiar with the basic functions of Artemis using a series of worked examples. These examples are designed to take you through the most immediately useful functions. However, there will be time, and encouragement, for you to explore other menus; features of Artemis that are not described in the exercises in this manual, but which may be of particular interest to some users. Like all the Modules in this workshop, please remember:

IF YOU DON' T UNDERSTAND, PLEASE ASK!

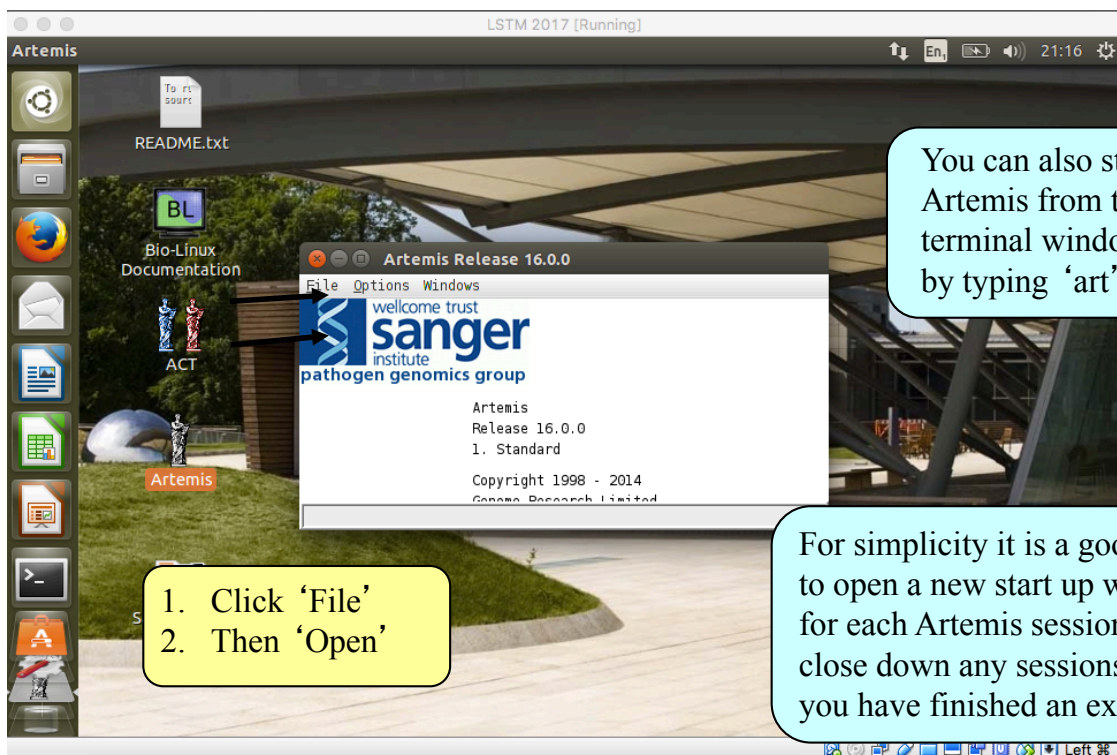
Artemis Exercise 1

1. Starting up the Artemis software

Double click the Artemis icon on the desktop.

A small start-up window will appear (see below). The directory **Module_1_Artemis** contains all files you will need for this module.

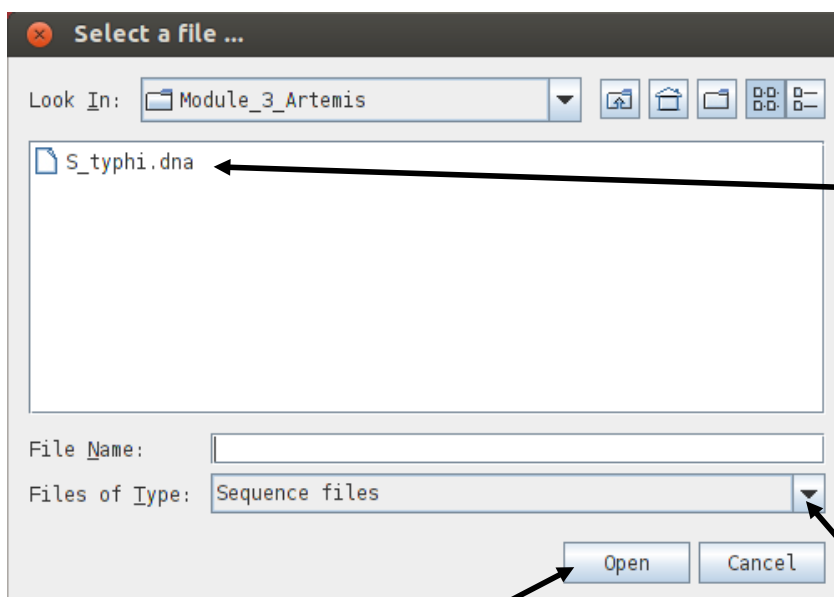
Now follow the sequence of numbers to load up the *Salmonella* Typhi chromosome sequence. Ask a demonstrator for help if you have any problems.



You can also start Artemis from the terminal window by typing 'art'

1. Click 'File'
2. Then 'Open'

For simplicity it is a good idea to open a new start up window for each Artemis session and close down any sessions once you have finished an exercise.



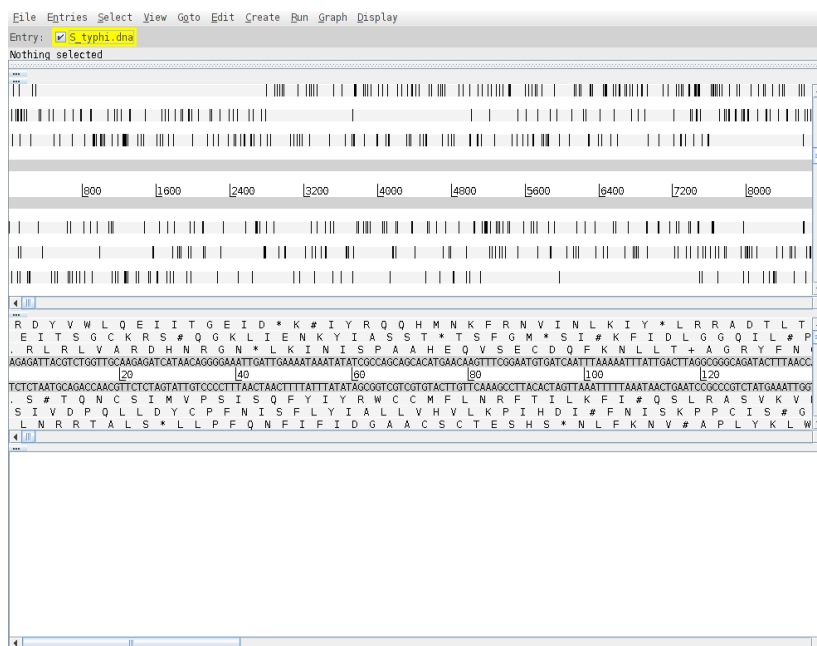
3
Single click to select file S_typhi.dna

Change to 'All Files' if you want to display all the files in the directory.
Use this feature to choose the type of file to be displayed in this panel. DNA sequence files will have the suffix '.dna'. Annotation files end with '.tab'. You can also open '.embl' files.

4 Single click to open file in Artemis then wait

2. Loading an annotation file (entry) into Artemis

Hopefully you will now have an Artemis window like this! If not, ask a demonstrator for assistance.



Now follow the numbers to load the annotation file for the *Salmonella* Typhi chromosome.

1

Click 'File' then 'Read an Entry'

What's an "Entry"?
It's a file of features which can be overlaid onto the sequence information displayed in the main Artemis view panel.

2

Single click to select file S_typhi.tab

3 Single click to open file in Artemis then wait (click 'no' if an error window pops up)

3. The basics of Artemis

Now you have an Artemis window open let's look at what is in there.

The screenshot displays the Artemis genome browser interface. At the top, a menu bar includes 'File', 'Entries', 'Select', 'View', 'Goto', 'Edit', 'Create', 'Run', 'Graph', and 'Display'. Below the menu bar, the 'Entry' section shows 'S_typhi.dna' and 'S_typhi.tab' with 'S_typhi.tab' selected. The 'Selected feature' section displays details for gene STY0004, including its location (bases 1287 to 428), amino acid count, and product (threonine synthase). The main sequence view panel shows the DNA sequence with forward and reverse strands, and three reading frames. Colored boxes represent various features, including CDSs and misc features. A zoomed-in view panel shows a detailed view of a CDS. The feature list panel at the bottom lists various features with their coordinates and descriptions. Numbered callouts (1-8) point to specific UI elements: 1. Menu bar, 2. Entry line, 3. Selected feature box, 4. Main sequence view panel, 5. Zoomed-in view panel, 6. Sliders for zooming, 7. Sliders for scrolling, 8. Slider for feature list.

1. **Drop-down menus:** There's lots in there so don't worry about all the details right now.
2. **Entry (top line):** shows which entries are currently loaded with the default entry highlighted in yellow (this is the entry into which newly created features are created). Selected feature: the details of a selected feature are shown here; in this case gene STY0004 (yellow box surrounded by thick black line).
3. This is the main **sequence view panel**. The central 2 grey lines represent the forward (top) and reverse (bottom) DNA strands. Above and below those are the 3 forward and 3 reverse reading frames. Stop codons are marked on the reading frames as black vertical bars. Genes and other annotated features (eg. Pfam and Prosite matches) are displayed as coloured boxes. We often refer to predicted genes as coding sequences or CDSs.
4. This panel has a similar layout to the main view panel but is zoomed in to show nucleotides and amino acids. Double click on a CDS in the main view to see the zoomed view of the start of that CDS. Note that both this and the main panel can be scrolled left and right (7, below) zoomed in and out (6, below).
5. **Feature panel:** This panel contains details of the various features, listed in the order that they occur on the DNA. Any selected features are highlighted. The list can be scrolled (8, below).
6. **Sliders** for zooming view panels.
7. **Sliders** for scrolling along the DNA.
8. **Slider** for scrolling feature list.

4. Getting around in Artemis

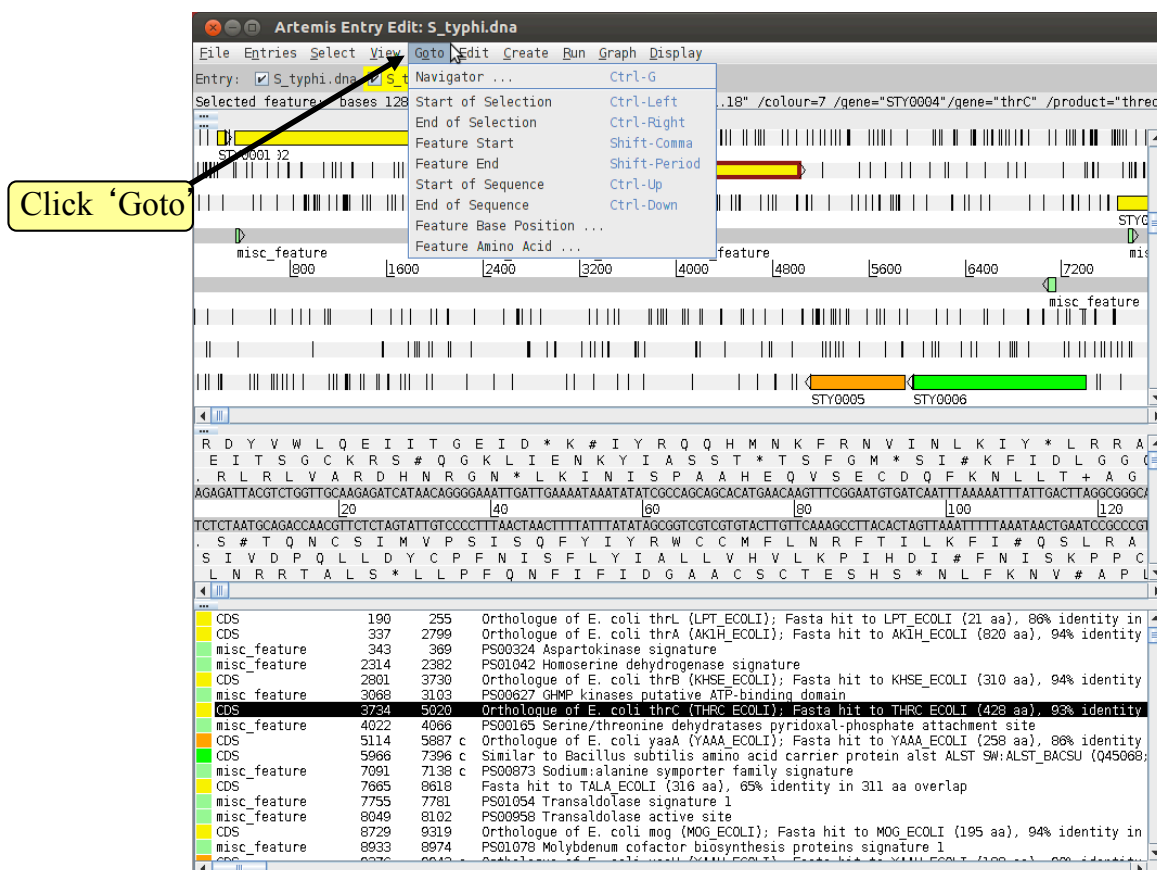
There are three main ways of getting to a particular DNA region in Artemis:

- the Goto drop-down menu
- the Navigator and
- the Feature Selector (which we will use in Exercise 4)

The best method depends on what you're trying to do. Knowing which one to use comes with practice.

4.1 The 'Goto' menu

The functions on this menu (below the Navigator option) are shortcuts for getting to locations within a selected feature or for jumping to the start or end of the DNA sequence. This is really intuitive so give it a try!



It may seem that 'Goto' 'Start of Selection' and 'Goto' 'Feature Start' do the same thing. Well they do if you have a feature selected but 'Goto' 'Start of Selection' will also work for a region which you have selected by click-dragging in the main window. So yes, give it a try!

Suggested tasks:

1. Zoom out, select / highlight a large region of sequence by clicking the left hand button and dragging the cursor then go to the start and end of this selected region.
2. Select a CDS then go to the start and end.
3. Go to the start and end of the genome sequence.
4. Select a CDS. Within it, go to a base (nucleotide) and/or amino acid of your choice.
5. Highlight a region then, from the right click menu, select 'Zoom to Selection'.

4.2 Navigator

The Navigator panel is fairly intuitive so open it up and give it a try.

Click 'Goto'
then Navigator

Check that the
appropriate search
button is on

The screenshot shows the Artemis software interface. The main window displays a genomic map with various features like CDS, misc_feature, and STY0001-0006. A 'Navigator' window is open, showing a list of search options: 'Goto Base:', 'Goto Feature With Gene Name:', 'Goto Feature With This Qualifier Value:', 'Goto Feature With This Key:', 'Find Base Pattern:', and 'Find Amino Acid String:'. The 'Goto Base:' option is selected. Below the search options, there are checkboxes for 'Start search at: beginning (or end) selection', 'Overlaps With Selection', 'Forward Strand', 'Reverse Strand', 'Search Backward', 'Ignore Case', and 'Allow Substring Matches'. The 'Goto' button is highlighted.

Suggestions about where to go:

1. Think of a number between 1 and 4809037 and go to that base (notice how the cursors on the horizontal sliders move with you).
2. Your favourite gene name (it may not be there so you could try '*fts*').
3. Use '**Goto Feature With This Qualifier value**' to search the contents of all qualifiers for a particular term. For example using the word 'pseudogene' will take you to the next feature with the word 'pseudogene' in any of its qualifiers. Note how repeated clicking of the 'Goto' button takes you to the following pseudogene in the order that they occur on the chromosome.
4. Look at **Appendix VI** which is a functional classification scheme used for the annotation of *S. Typhi*. Each CDS has a class qualifier best describing its function. Use the '**Goto Feature With This Qualifier value**' search to look for CDSs belonging to a class of interest by searching with the appropriate class values.
5. tRNA genes. Type 'tRNA' in the '**Goto Feature With This Key**'.
6. Regulator-binding DNA consensus sequence (real or made up!). Note that degenerate base values can be used (**Appendix VIII**).
7. Amino acid consensus sequences (real or made up!). You can use 'x's. Note that it searches all six reading frames regardless of whether the amino acids are encoded or not.

What are Keys and Qualifiers? See **Appendix IV**

Clearly there are many more features of Artemis which we will not have time to explain in detail. Before getting on with this next section it might be worth browsing the menus. Hopefully you will find most of them easy to understand.

Artemis Exercise 2

This part of the exercise uses the files and data you already have loaded into Artemis from Part I. By a method of your choice go to the region from bases 2188349 to 2199512 on the DNA sequence. This region is bordered by the *fbxB* gene which codes for fructose-bisphosphate aldolase. You can use the Navigator function discussed previously to get there. The region you arrive at should look similar to that shown below.

The screenshot shows the Artemis interface with the following components:

- Top Panel:** File Entries Select View Goto Edit Create Run Graph Display. Entry: S_typhi.dna S_typhi.tab. Nothing selected.
- Genomic Map:** A horizontal bar representing the DNA sequence with various features. Two callout boxes on the right point to specific features:
 - CDS features:** Points to green bars representing coding sequences (e.g., STY2343, STY2345, STY2348, STY2349, STY2365, STY2371, STY2373).
 - Misc features:** Points to orange and yellow bars representing other features (e.g., RBS, misc_feature).
- DNA Sequence:** A text view of the DNA sequence with a green highlight under the sequence:


```

            N # Y L Y N N K P I V # T A W N Q Q E E P I P Y S F D Y Y N # H L L K # C Q T I W P N A V N
            I N I F I I T S Q L C K L R G I N R K N Q F L T L L I I T I N T Y # N N A K Q Y G R T Q # T
            E L I S L # # Q A N C V N C V E S T G R T N S L L F * L L Q L T L I K I M P N N M A E R S K I
            A A T T A A T A T C T T T A T A A T A A C A A G C C A A T T G T G T A A A C T G C G T G G A A T C A C A G G A A G A A C C A A T T C C T T A C T C T T T G A T T A T T A C A A T T A A C A C T T A T T A A A A T A A T G C C A A C A A T A T G C C G A A C G C A G T A A A C
            [2188940] [2188960] [2188980] [2189000] [2189020] [2189040] [2189060]
            T T A A T T A T A G A A A T A T T G T T C G G T T A A C A C A T T G A C G C A C C T A G T T G T C C T T C T G G T T A A G G A A T G A G A A A C T A A T A A T G T A A T T G T G A A T A T T T A T T A C G G T T G T T A T A C C G C C T T G C G T C A T T T G
            F # Y R # L L L L G I T Y V A H F * C S S G I G # E K S # # L # C K N F Y H W V I H G F A T F I
            I L I K I I V L W N H L S R P I L L F F W N R V R K I I V I L V # F L A L C Y P R V C Y V
            N I D K Y Y C A L Q T F O T S D V P L V L W L K S K O N N C N V S I L I T G F L I A S R L L C
            
```
- Feature Table:** A table at the bottom listing features with their coordinates and descriptions:

CDS	190	255	Orthologue of E. coli thrL (LPT_ECOLI); Fasta hit to LPT_ECOLI (21 aa), 86% identity in 21 aa overlap
CDS	337	2799	Orthologue of E. coli thrA (AKIH_ECOLI); Fasta hit to AKIH_ECOLI (820 aa), 94% identity in 820 aa overlap
misc_feature	343	369	PS00324 Aspartokinase signature
misc_feature	2314	2382	PS01042 Homoserine dehydrogenase signature
CDS	2801	3730	Orthologue of E. coli thrB (KHSE_ECOLI); Fasta hit to KHSE_ECOLI (310 aa), 94% identity in 308 aa overlap
misc_feature	3068	3103	PS00627 GMP kinases putative ATP-binding domain
CDS	3734	5020	Orthologue of E. coli thrC (THRC_ECOLI); Fasta hit to THRC_ECOLI (428 aa), 93% identity in 428 aa overlap
misc_feature	4022	4066	PS00165 Serine/threonine dehydratases pyridoxal-phosphate attachment site
CDS	5114	5887	Orthologue of E. coli yaaA (YAAA_ECOLI); Fasta hit to YAAA_ECOLI (258 aa), 86% identity in 257 aa overlap
CDS	5966	7396	Similar to Bacillus subtilis amino acid carrier protein alst ALST SW:ALST_BACSU (045068; P40743) fasta hit to ALST_BACSU (1450 aa), 86% identity in 1450 aa overlap
misc_feature	7091	7138	PS00873 Sodium:alanine symporter family signature
CDS	7665	8618	Fasta hit to TALA_ECOLI (316 aa), 65% identity in 311 aa overlap
misc_feature	7755	7781	PS01054 Transaldolase signature 1
misc_feature	8049	8102	PS00958 Transaldolase active site
CDS	8729	9319	Orthologue of E. coli mog (MOG_ECOLI); Fasta hit to MOG_ECOLI (195 aa), 94% identity in 192 aa overlap
misc_feature	8933	8974	PS01078 Molibdenum cofactor biosynthesis proteins signature 1

Once you have found this region have a look at some of the information available:

Information to view:

Annotation

If you click on a particular feature you can view the annotation associated with it: select a CDS feature (or any other feature) and click on the 'Edit' menu and select 'Selected Feature in Editor'. A window will appear containing all the annotation that is associated with that CDS. The format for this information is constrained by that which can be submitted to the EMBL database.

Viewing amino acid or protein sequence

Click on the 'View' menu and you will see various options for viewing the bases or amino acids of the feature you have selected, in two formats i.e. EMBL or fasta. This can be very useful when using other programs that are not integrated into Artemis e.g. those available on the Web that require you to cut and paste sequence into them.

Plots/Graphs

Feature plots can be displayed by selecting a CDS feature then clicking 'View' and 'Feature Plots'. The window which appears shows plots predicting hydrophobicity, hydrophilicity and coiled-coil regions for the protein product of the selected CDS.

Load additional files

You should be able to see the results from Prosite searches, run on the translation of each CDS, as pale-green boxes on the grey DNA lines. The results from the Pfam protein motif searches are not yet shown, but can be viewed by loading the appropriate file. Click on 'File' then 'Read an Entry' and select the file PF.tab. Each Pfam match will appear as a coloured blue feature in the main display panel on the grey DNA lines. To see the details click the feature then click 'View' then 'Selection' or click 'Edit' then 'Selected Features in Editor'. Please ask if you are unsure about Prosite and Pfam.

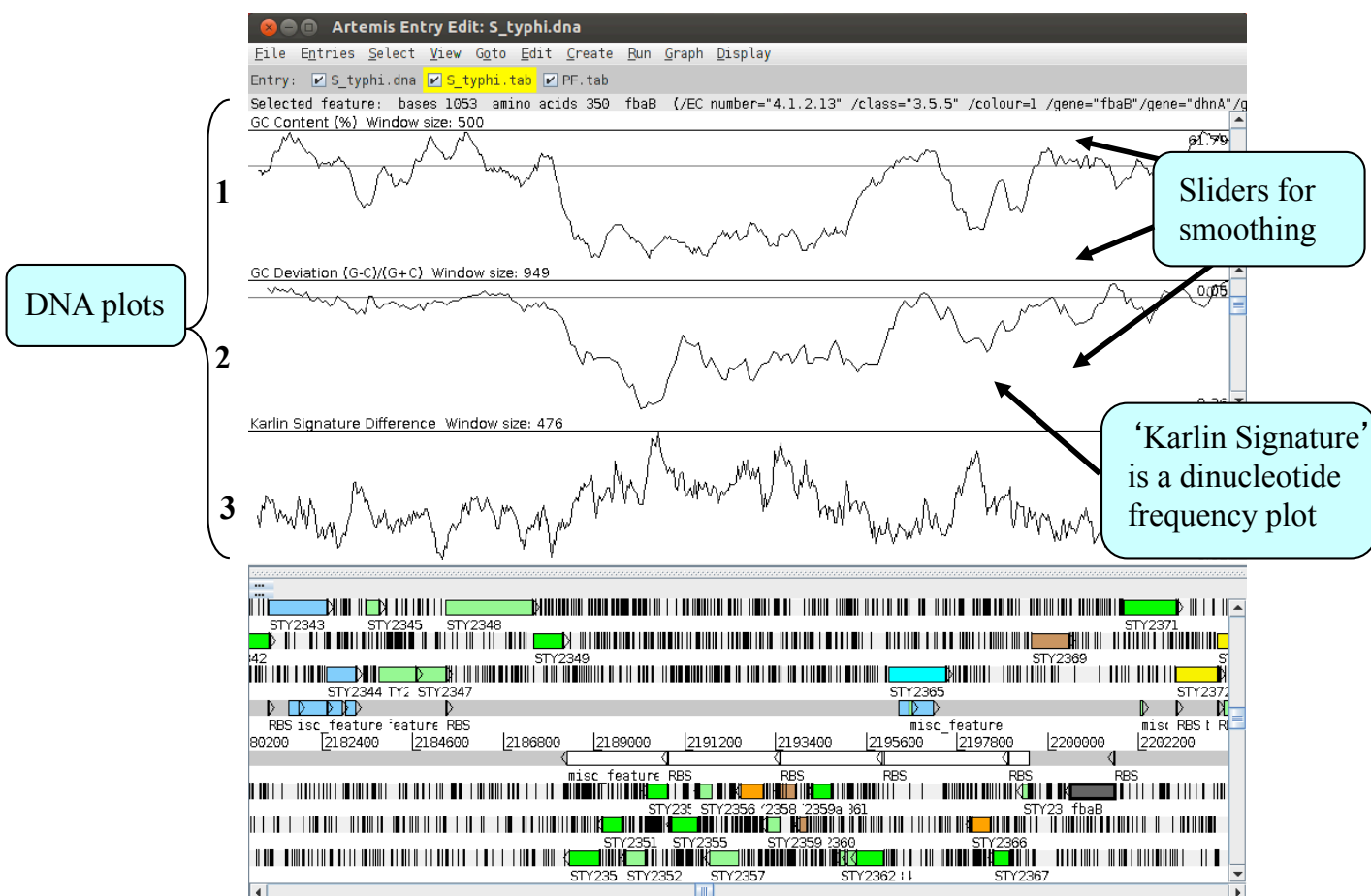
Further information on specific Prosite or Pfam entries can be found on the web at:
<http://ca.expasy.org/prosite> and <http://pfam.sanger.ac.uk/>

In addition to looking at the fine detail of the annotated features it is also possible to look at the characteristics of the DNA covering the region displayed. This can be done by adding various plots to the display, showing different characteristics of the DNA. Some of the plots can be used to look at the protein coding potential of translation frames within the DNA, such as GC frame plot, and others can be used to search for horizontally acquired DNA.

The plot information is generated dynamically by Artemis and although this is a relatively speedy exercise for a small region of DNA, on a whole genome view (we will move onto this later) this may take a little time, so be patient.

To view the graphs:

Click on the 'Graph' menu to see all those available. Perhaps some of the most useful plots are the (1) 'GC Content (%)', (2) 'GC Deviation' and (3) 'Karlín Signature Difference' as shown below. To adjust the smoothing of the graph you change the window size over which the points on the graph are calculated, using the sliders shown below. If you are not familiar with any of these please ask.



Notice how several of the plots show a marked deviation around the region you are currently looking at. To fully appreciate how anomalous this region is move the genome view by scrolling to the left and right of this region. The apparent unusual nucleotide content of this region is indicative of laterally acquired DNA that has inserted into the genome.

Your Artemis window should now look similar to the one shown.

As well as looking at the characteristics of small regions of the genome, it is possible to zoom out and look at the characteristics of the genome as a whole. To view the entire genome you can use the sliders indicated below. However, be careful zooming out quickly with all the features being displayed, as this may temporarily lock up the computer.

1. To make this process faster and clearer, **switch off stop codons** by clicking with the right mouse button in the main view panel. A menu will appear with an option to de-select 'Stop Codons' (see below).

2. You will also need to temporarily **remove all of the annotated features** from the Artemis display window. In fact if you leave them on, which you can, they would be too small to see when you zoomed out to display the entire genome. To remove the annotation click on the S_typhi.tab entry button on the grey entry line of the Artemis window shown above.

2 To de-select the annotation click here.

No stop codons shown on frame lines

1 Menu item for de-selecting stop codons

Artemis Entry Edit: S_typhi.dna

File Entries Select View Goto Edit Create Run Graph Display

Entry: S_typhi.dna S_typhi.tab PF.tab

Selected feature: bases 1053 amino acids 350 fbaB (/EC number="4.1.2.13" /class="3.5.5" /colour=1 /gene="fbaB"/gene="dhxA"/g

GC Content (%) Window size: 500

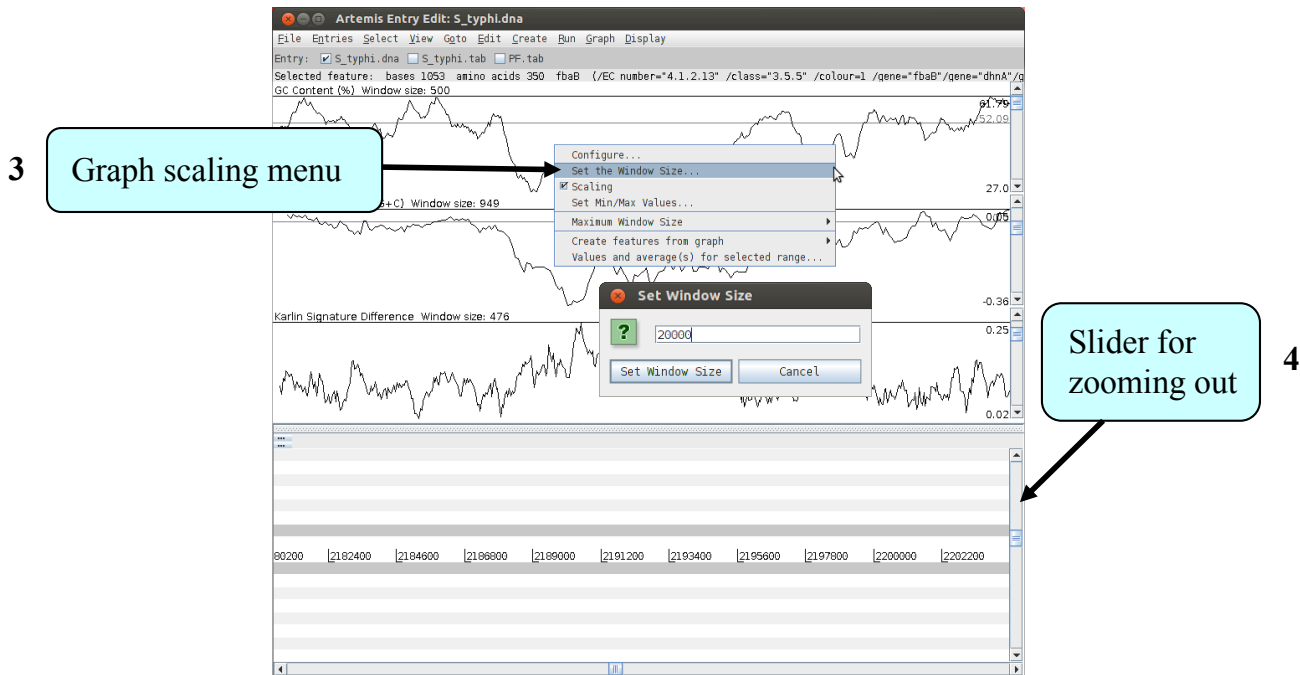
GC Deviation (G-C)/(G+C) Window size: 949

Karlin Signature Difference Window size: 476

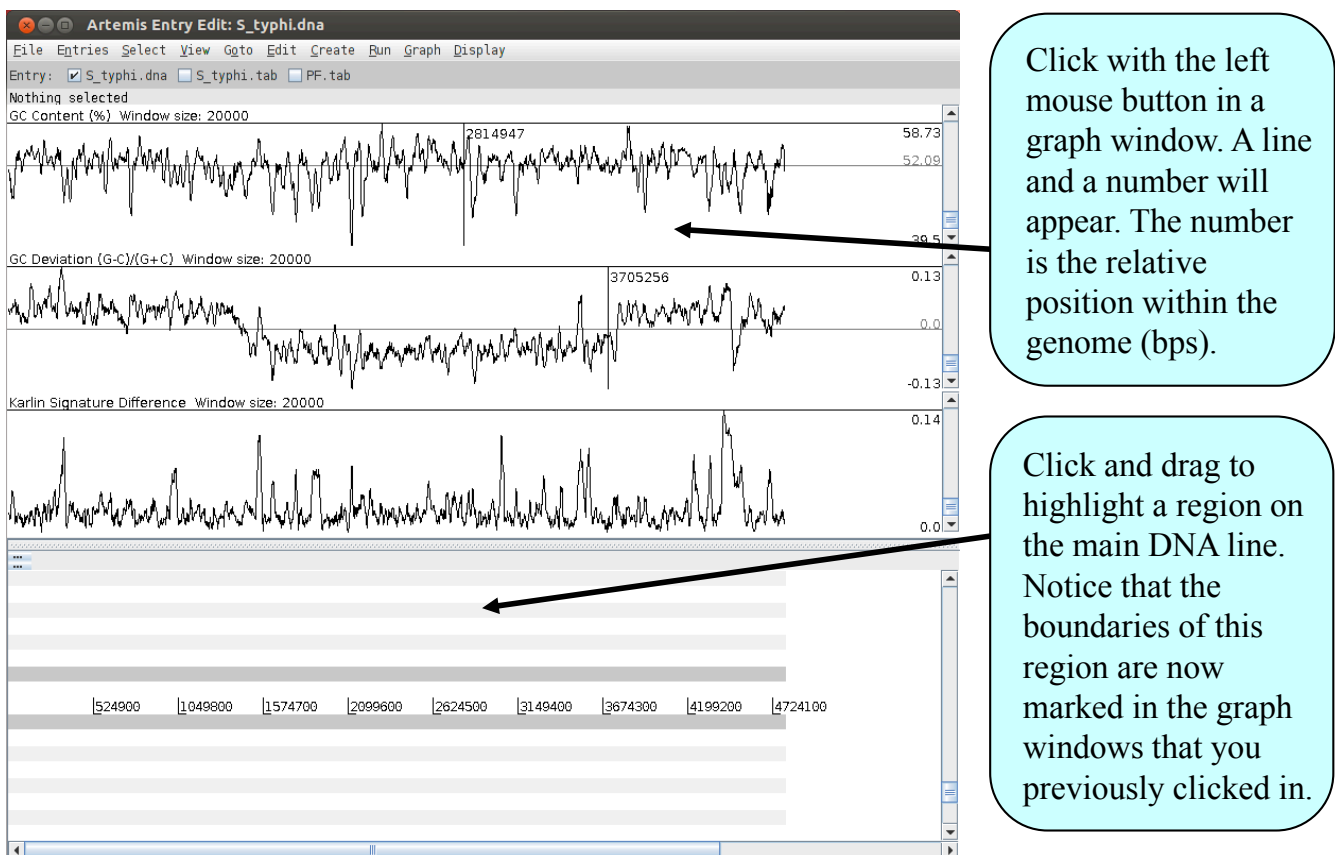
Smallest Features In Front
Set Score Cutoffs ...
Raise Selected Features
Lower Selected Features
Zoom to Selection
Select Visible Range
Select Visible Features
Frame Line Features ...

Entries
Select
View
Goto
Edit
Create
Write
Run

Start Codons
 Stop Codons
 Feature Arrows
 Feature Borders
 Feature Labels
 One Line Per Entry
 Feature Stack View
 Forward Frame Lines
 Reverse Frame Lines
 All Features On Frame Lines
 Show Source Features
 Flip Display
 Colourise Bases

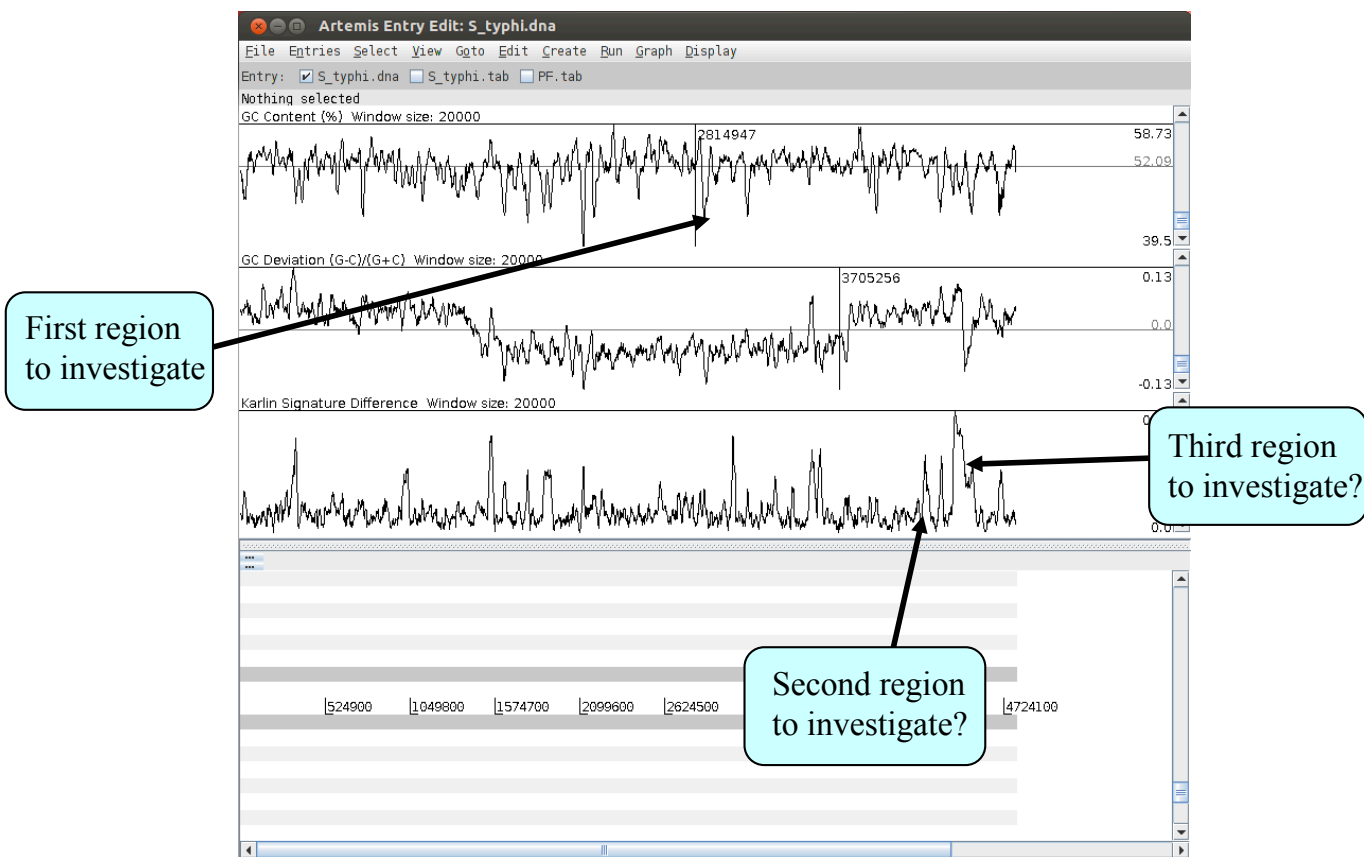


3. One final tip is to **adjust the scaling** for each graph displayed before zooming out. This increases the maximum window size over which a single point for each plot is calculated. To adjust the scaling click with the right mouse button over a particular graph window. A menu will appear with an option “Set the Window size’ (see above), set the window size to ‘20000’. You should do this for each graph displayed (if you get an error message press continue).
4. You are now ready to zoom out by dragging or clicking the slider indicated above. Once you have zoomed out fully to see the entire genome you will need to adjust the smoothing of the graphs using the vertical graph sliders as before, to have a similar view to that shown below.



Artemis Exercise 3

There are many examples where anomalous regions of DNA within a genome have been shown to carry laterally acquired DNA. In this part of the exercise we are going to look at several of these regions in more detail. Starting with the whole genome view, note down the approximate positions and characteristics of the three regions indicated above. Remember the locations of the peaks are given in the graph window if you click the left mouse button within it.

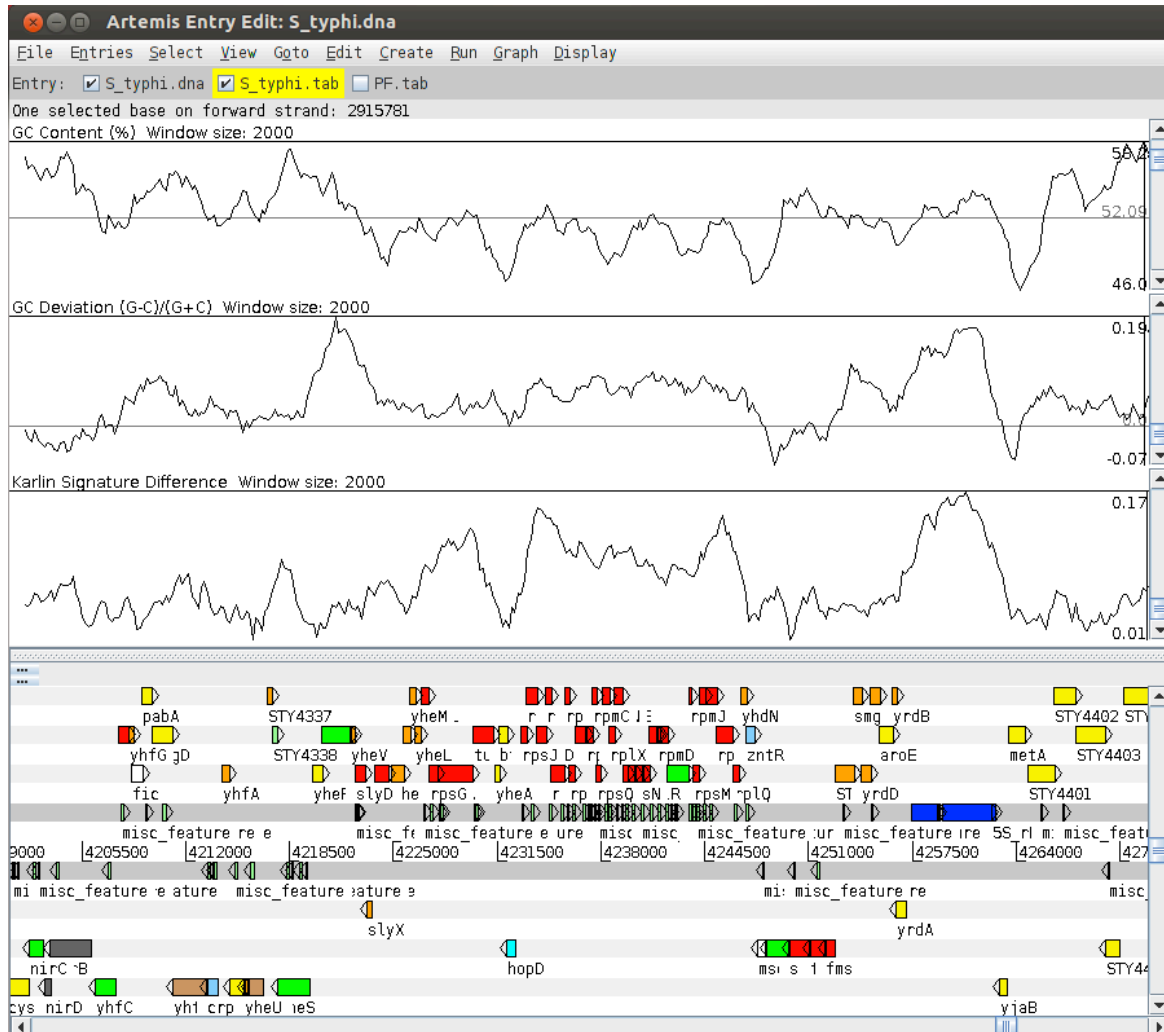


Genome location	Characteristics of DNA plots
Region 1 : 2,860,000 bps	peak - karlin, troughs for G+C and CG deviation
Region 2 :	
Region 3 :	

We will now zoom back into the genome to look in more detail at the first of these three peaks. Using the left mouse button, highlight the anomalous region of the graph - this will also highlight the region in the main display. You can then use the 'right mouse button menu' in the main display to 'Zoom to selection' - you may need to zoom out from there. Remember that in order to see the CDS features lying within this region you will need to turn the annotation (S_typhi.tab) entry back on.

Use one of the methods you have already used to take you to the second region of interest that you noted down.

Region 2



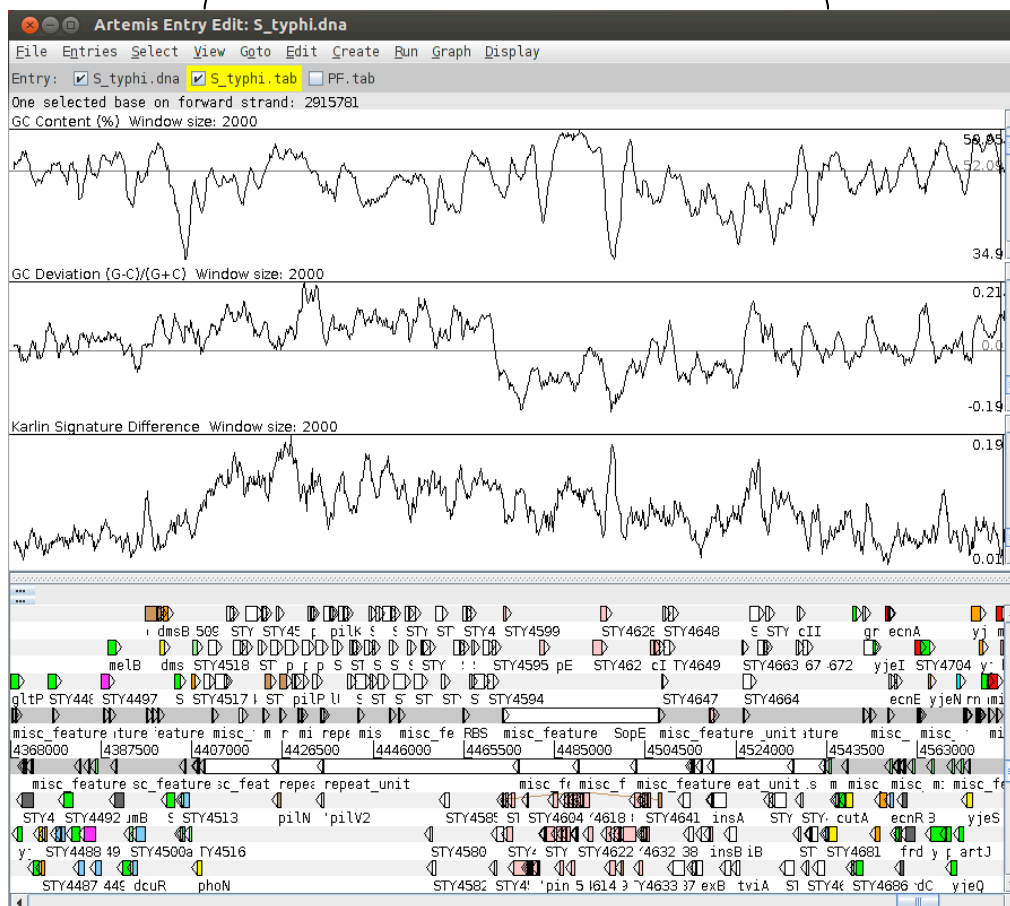
Region two acts as a cautionary note when looking at anomalous regions within a genome. Have a look at the features and annotation of the CDSs within this region:

- Does this region have any of the characteristics of pathogenicity island?
- Are the genes within this region essential or dispensable (“accessory”)?

Is it possible that the atypical base composition of this region is not a consequence of having originated from a foreign host? The base composition may actually be reflective of the tight sequence constraints under which this region has been maintained, in contrast to the background level sequence variation in the rest of the genome.

Next go to Region 3.

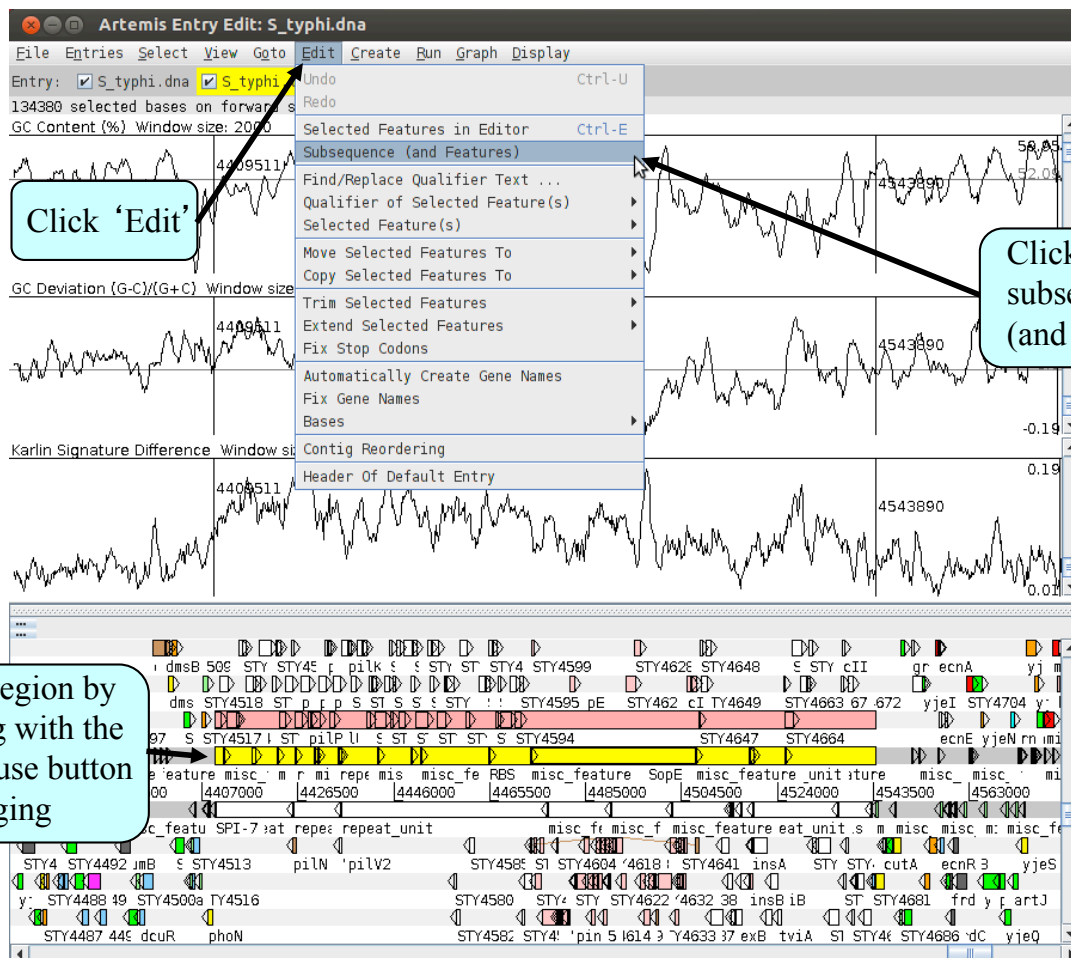
Region 3



As with region 1, this region is also defined as a *Salmonella* pathogenicity island (SPI). SPI-7, or the major Vi pathogenicity island, is ~134 kb in length and contains ~30 kb of integrated bacteriophage. Have a look at the CDSs within this region. As before notice any stable RNAs that may have acted as the phage integration site.

Artemis Exercise 4

Continuing on from the analysis of Region 3 or SPI-7 (the major Vi-antigen pathogenicity island) we are going to extract this region from the whole genome sequence and perform some more detailed analysis on it. We will aim to write and save new EMBL format files which will include just the annotations and DNA for this region. Follow the numbers on the next page to complete the task.

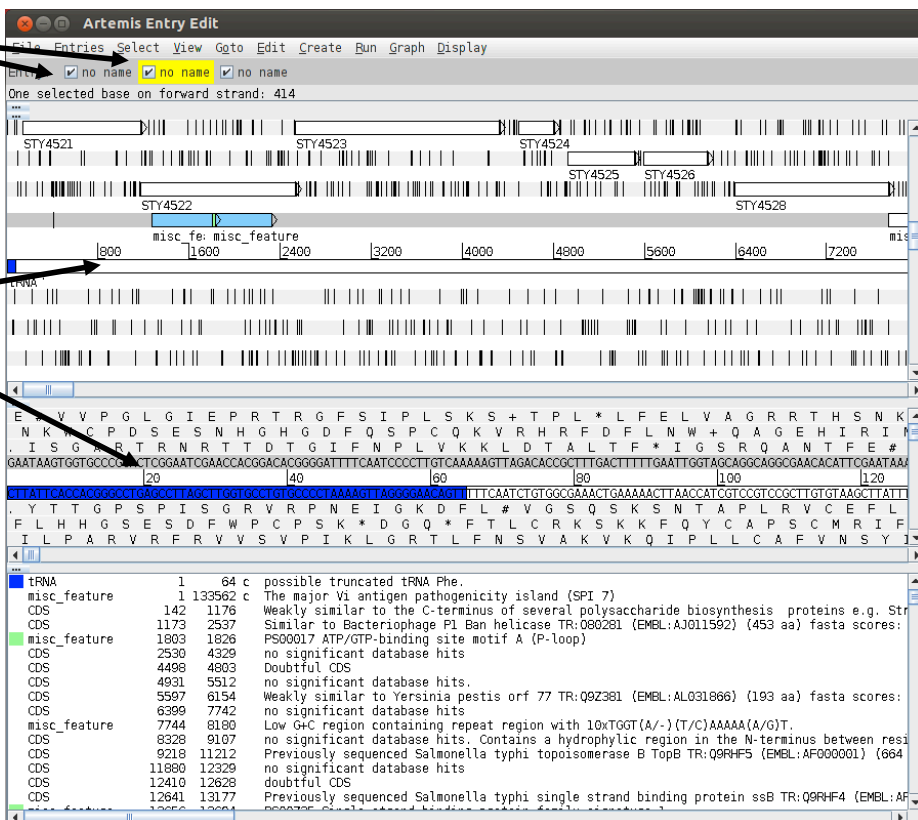


1 Select region by clicking with the left mouse button & dragging

A new Artemis window will appear displaying only the region that you highlighted

Note the entry names have changed

Note the bases have been renumbered from the first base you selected.



Note that the two entries on the grey 'Entry' line are now denoted 'no name'. They represent the same information in the same order as the original Artemis window but simply have no assigned 'Entry' names. As the sub-sequence is now viewed in a new Artemis session, this prevents the original files (S_typhi.dna and S_typhi.tab) from being over-written.

We will save the new files with relevant names to avoid confusion. So click on the 'File' menu then 'Save An Entry As' and then 'New File'. Another menu will ask you to choose one of the entries listed. At this point they will both be called 'no name'. Left click on the top entry in the list. A window will appear asking you to give this file a name. Save this file as spi7.dna

Do the same again for the second unnamed entry and save it as spi7.tab

Feature	Start	End	Description
CDS	1	548	Similar to Salmonella typhimurium nonspecific acid phosphatase precursor phoN SW:PHON_SALTY (P26)
misc_feature	72	95	PS01157 Class A bacterial acid phosphatases signature
tRNA	1252	1315	possible truncated tRNA Phe.
misc_feature	1252	134813	The major Vi antigen pathogenicity island (SPI 7)
CDS	1393	2427	Weakly similar to the C-terminus of several polysaccharide biosynthesis proteins e.g. Streptococcus
CDS	2424	3798	Similar to Bacteriophage P1 Ban helicase TR:080281 (EMBL:AJ011592) (453 aa) fasta scores: E(): 0.
misc_feature	3054	3077	PS00017 ATP/GTP-binding site motif A (P-loop)
CDS	3781	5590	no significant database hits
CDS	5749	6054	Doubtful CDS
CDS	6182	6763	no significant database hits.
CDS	6848	7405	Weakly similar to Yersinia pestis orf 77 TR:Q9Z381 (EMBL:AL031866) (193 aa) fasta scores: E(): 8
CDS	7650	8993	no significant database hits
misc_feature	8995	9431	Low G+C region containing repeat region with 10xTGGT(A/-)(T/C)AAAA(A/G)T.
CDS	9579	10358	no significant database hits. Contains a hydrophilic region in the N-terminus between residues 34
CDS	10469	12463	Previously sequenced Salmonella typhi topoisomerase B TopB TR:Q9RHF5 (EMBL:AF000001) (664 aa) fas
CDS	13131	13580	no significant database hits
CDS	13661	13879	doubtful CDS
CDS	13892	14428	Previously sequenced Salmonella typhi single strand binding protein ssB TR:Q9RHF4 (EMBL:AF000001)

We are going to look at this region in more detail and to attempt to define the limits of the bacteriophage that lies within this region. Luckily for us all the phage-related genes within this region have been given a colour code number 12 (pink; for a list of the other numerical values that Artemis will display as colours for features see **Appendix VII**). We are going to use this information to select all the relevant phage genes using the Feature selector as shown below and then define the limits of the bacteriophage.

First we need to create a new entry (click 'Create' then 'New Entry'). Another entry will appear on the entry line called, you guessed it, 'no name'. We will eventually copy all our phage-related genes into here.

The image shows a sequence of six numbered steps for using the Artemis Feature Selector:

- Click 'Select' then 'Feature Selector'**: The 'Select' menu is highlighted in the top bar, and the 'Feature Selector...' option is chosen from the dropdown.
- Make sure the buttons are selected**: The 'Feature Selector' dialog box is open, and the 'Key' and 'Qualifier' fields are selected.
- Set Key to 'CDS' and Qualifier to 'colour'**: The 'Key' dropdown is set to 'CDS' and the 'Qualifier' dropdown is set to 'colour'.
- Type search term**: The 'Containing this text' field is filled with the number '12'.
- Click to select features containing search term**: The 'Select' button at the bottom of the dialog is clicked.
- Click to view selected features in a list**: The 'View' button at the bottom of the dialog is clicked.

The final window, titled 'All features with key "CDS" with qualifier "colour" containing text "12"', displays a list of features:

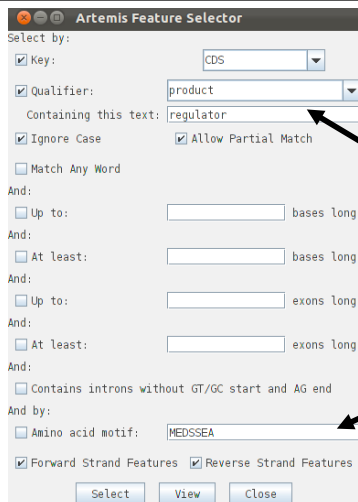
Accession	Key	Qualifier	Description
65714	CDS		no significant database hits
66178	CDS		Similar to Escherichia coli prophage P2 Ogr protein SW:OGRK.E
66464	CDS		Similar to Bacteriophage P2 late gene control protein D SW:VPI
67561	CDS		Similar to Bacteriophage P2 complete genome U essential tail p
69046	CDS		Similar to Bacteriophage 186 protein G TR:Q37049 (EMBL:U32222)
70810	CDS		Similar to Bacteriophage 186 Orf52 H TR:080316 (EMBL:U32222)
70953	CDS		Similar to Bacteriophage P2 complete genome E, essential tail
71310	CDS		Similar to Bacteriophage P2 major tail tube protein fII SW:VPI
71835	CDS		Similar to Bacteriophage P2 major tail sheath protein fI SW:VPI
73542	CDS		Similar to Salmonella typhimurium invasion-associated secreted
74462	CDS		Similar to Bacteriophage P2 probable tail fiber assembly prot
74876	CDS		Similar to Bacteriophage P2 probable tail fiber protein SW:VPI
76492	CDS		Similar to Bacteriophage P2 tail protein I SW:VPI_BPP2 (P2670)
77090	CDS		Similar to Bacteriophage P2 baseplate assembly protein J SW:VPI
77985	CDS		Similar to Bacteriophage P2 baseplate assembly protein W SW:VPI
78341	CDS		Similar to Bacteriophage P2 baseplate assembly protein V SW:VPI
78968	CDS		Similar to Bacteriophage P2 tail completion protein S SW:VPI
79427	CDS		Similar to Bacteriophage P2 tail completion protein R SW:VPI
79594	CDS		Similar to Bacteriophage P2 protein LysB protein involved in
80379	CDS		no significant database hits, contains possible membrane span
80761	CDS		Similar to Serratia marcescens putative phage lysozyme NucD TH

The genes listed in (6) are only those fitting your selection criteria. They can be copied or cut / moved in to a new entry so we can view them in isolation from the rest of the information within spi7.tab.

Firstly in window (6) select all of the CDSs shown by clicking on the 'Select' menu and then selecting 'All'. All the features listed in window (6) should now be highlighted. To copy them to another entry (file) click 'Edit' then 'Copy Selected Features To' then 'no name'. Close the two smaller feature selector windows and return to the SPI-7 Artemis window. You could rename the 'no name' entry as phage.tab, as you did before. Temporarily remove the features contained in 'spi7.tab' file by left clicking on the entry button on the grey entry line. Only the phage genes should remain.

Additional methods for selecting/extracting features using the Feature Selector

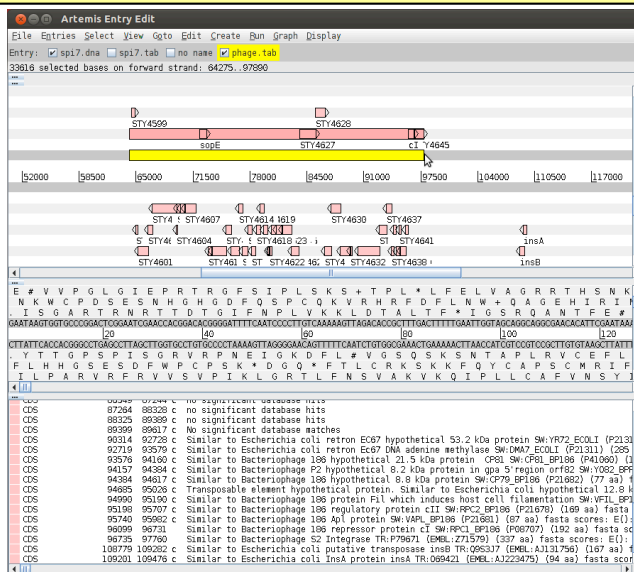
It is worth noting that the Feature Selector can be used in many other ways to select and extract subsets of features from the genome, using eg text or amino acid searches.



Space for a search term or amino acid motif

Defining the extent of the prophage

Even from this preliminary analysis it is clear that the prophage occupies a fairly discrete region within SPI-7 (see below). It is often useful to create a new DNA feature to define the limits of this type of genome landmark. To do this use the left mouse button to click and drag over the region that you think defines the prophage.



While the region is highlighted, click on the 'Create' menu and select 'Create feature from base range'. A feature edit window will appear. The default 'Key' value given by Artemis when creating a new feature is 'CDS'. With this 'Key' the newly created feature would automatically be put on the translation line. However, if we change this to 'misc_feature' (an option in the 'Key' drop down menu in the top left hand corner of the Edit window), Artemis will place this feature on the DNA line. This is perhaps more appropriate and is easier to visualise. You can also add a qualifier, such as '/label': select 'label' from the 'Add Qualifier' list and click 'Add Qualifier', '/label=' will appear in the text window; add text of your choice, then click 'OK'. That text will be used as a feature label to be displayed in the main sequence view panel.

To see how well you have done, turn the spi7.tab.

Your final task is to write out the spi7 files in EMBL submission format, and create a merged annotation and sequence file in EMBL submission format. In Artemis you are going to copy the annotation features from the ‘.tab’ file into the ‘.dna’ file, and then save this entry in EMBL format. Don’t worry about error messages popping up. This is because not all entries are accepted by the EMBL database.

1 Click 'Select' then 'All'

2 Click 'Edit', then 'Copy Selected Features To'

3 Select 'spi7.dna'

4 Click 'File' then 'Save An Entry As'

5 'EMBL Submission Format'

6 Select 'spi7.dna'

7 Save file as spi7.embl

Now open the EMBL format file that you have just created in Artemis.

tRNA	1	64	c	possible truncated tRNA Phe.
misc_feature	1	133562	c	The major Vi antigen pathogenicity island (SPI 7)
CDS	142	1176		Weakly similar to the C-terminus of several polysaccharide biosynthesis proteins e.g. Str
CDS	1173	2537		Similar to Bacteriophage P1 Ban helicase TR:080281 (EMBL:AJ011592) (453 aa) fasta scores:
misc_feature	1803	1826		PS00017 ATP/GTP-binding site motif A (P-loop)
CDS	2530	4329		no significant database hits
CDS	4498	4803		Doubtful CDS
CDS	4931	5512		no significant database hits.
CDS	5597	6154		Weakly similar to Yersinia pestis orf 77 TR:Q9Z381 (EMBL:AL031866) (193 aa) fasta scores:
CDS	6399	7742		no significant database hits
misc_feature	7744	8180		Low G+C region containing repeat region with 10xTGGT(A/-)(T/C)AAAAA(A/G)T.
CDS	8328	9107		no significant database hits. Contains a hydrophilic region in the N-terminus between resi
CDS	9218	11212		Previously sequenced Salmonella typhi topoisomerase B TopB TR:Q9RHF5 (EMBL:AF000001) (664
CDS	11880	12329		no significant database hits
CDS	12410	12628		doubtful CDS
CDS	12641	13177		Previously sequenced Salmonella typhi single strand binding protein ssB TR:Q9RHE4 (EMBL:AF

You will see that the colours of the features have now changed. This is because not all the qualifiers in the previous entry are accepted by the EMBL database, so some have not been saved in this format. This includes the '/colour' qualifier, so Artemis displays the features with default colours.

When you download sequence files from EMBL and visualize them in Artemis you will notice that they are displayed using default colours. You can customize your own annotation files with the '/colour' qualifier and chosen number (**Appendix VII**), to differentiate features. To do this you can use the Feature Selector to select certain features and annotate them all using the 'Edit', 'Change Qualifiers of Selected' function.

Artemis Exercise 5

This exercise will introduce you to database searches and will give you a first insight in the annotation of genes.

The gene you will work on is *hpcC* (STY1136). Go to this gene by using one of the different methods you have learned so far. You will need to close down the last Artemis exercise if you haven't already done so. Start a new Artemis Session, as before, and open again *S_typhi.dna* and read the annotation file (*S_typhi.tab*) in.

As you can see the gene is full of stop codons indicating that we are looking at a pseudogene. To correct the annotation we are going to use database search. Follow now the numbers in the figure below to start a database search. The search may take a couple of minutes to run; a banner will pop up to tell you when its complete (3).

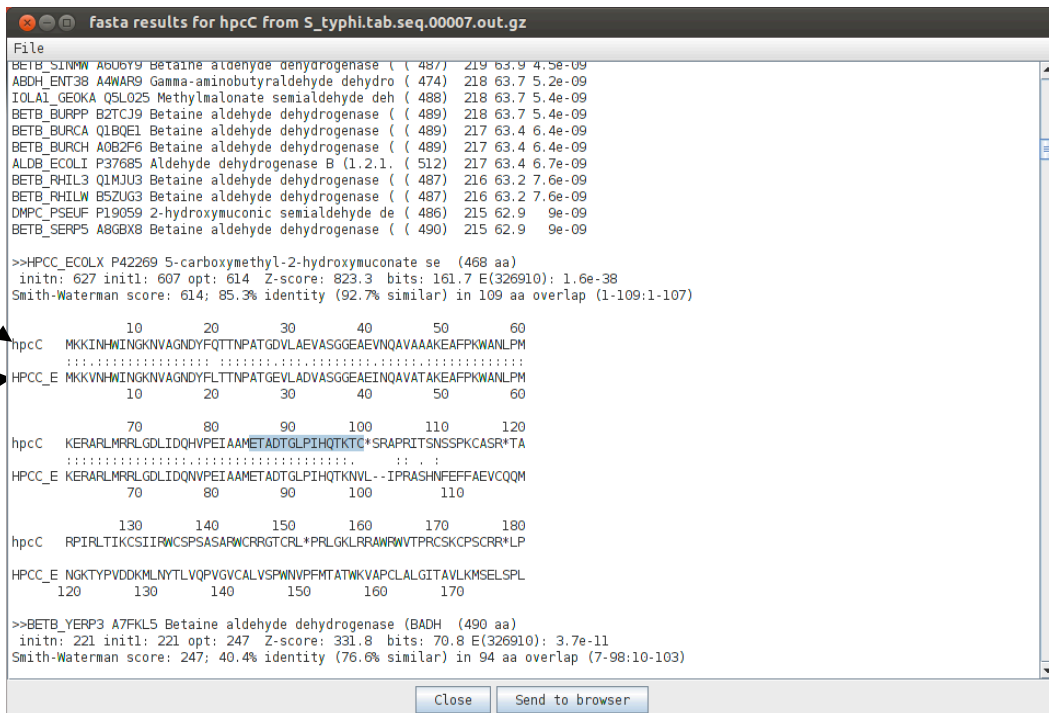
1 Select CDS

2 Start fasta against the Uniprot bacteria database

3

Feature Type	Start	End	Description
CDS	1099596	1101061	Pseudogene. Similar to Escherichia coli 5-carboxymethyl-2-hydroxyruconate semialdehyde d
misc_feature	1100321	1100344	PS00687 Aldehyde dehydrogenases glutamic acid active site
misc_feature	1100405	1100440	PS00670 Aldehyde dehydrogenases cysteine active site
CDS	1101063	1101914	Similar to Escherichia coli 3,4-dihydroxyphenylacetate 2,3-dioxygenase hpcB SW:HPCB_ECOL
CDS	1101924	1102304	Similar to Escherichia coli 5-carboxymethyl-2-hydroxyruconate delta-isomerase hpcD SW:HP
CDS	1102448	1103251	Similar to Escherichia coli 2-oxo-hepta-3-ene-1,7-dioic acid hydratase hpcG SW:HPCG_ECOL
CDS	1103262	1104053	Similar to Escherichia coli 2,4-dihydroxyhept-2-ene-1,7-dioic acid aldolase hpcH or hpaI
CDS	1104125	1105501	Similar to Escherichia coli putative 4-hydroxyphenylacetate permease hpaX TR:Q46984 (EMBL
CDS	1105511	1106407	Similar to Escherichia coli 4-hydroxyphenylacetate 3-monooxygenase operon regulatory pro
misc_feature	1106237	1106365	PS00041 Bacterial regulatory proteins, arnC family signature
CDS	1106421	1107359	No significant database matches
CDS	1108783	1109088	Orthologue of E. coli yccD (YCCD_ECOLI); Fasta hit to YCCD_ECOLI (101 aa), 74% identity
CDS	1109088	1110008	Fasta hit to DNAJ_ECOLI (375 aa), 35% identity in 353 aa overlap
misc_feature	1109814	1109873	PS00636 Nt-dnaJ domain signature
CDS	1110244	1110606	Similar to Salmonella typhimurium suppressor for copper-sensitivity a ScaA TR:033917 (EMBL
CDS	1110655	1112541	Similar to Salmonella typhimurium suppressor for copper-sensitivity B precursor ScaB TR:

To view the search results click 'View', then 'Search Results', then 'fasta results'. The results will appear in a scrollable window. Scroll down to the first sequence comparison and you should see the results as shown in the next figure.

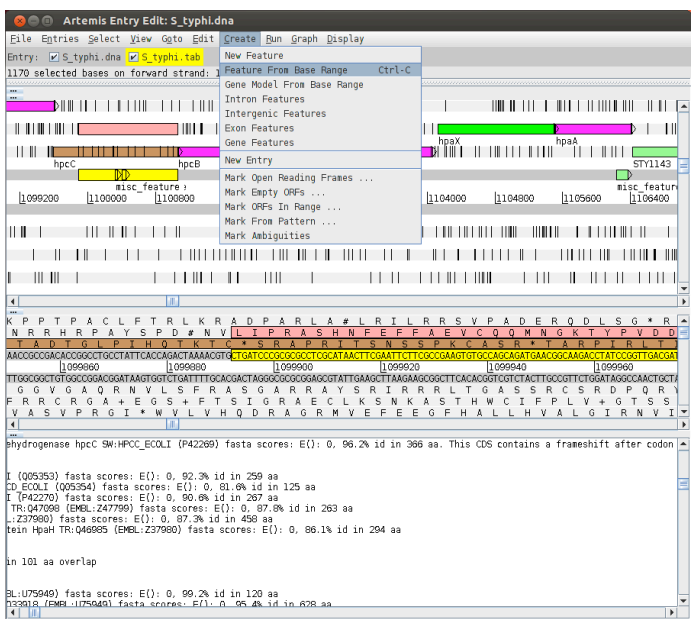


Our gene

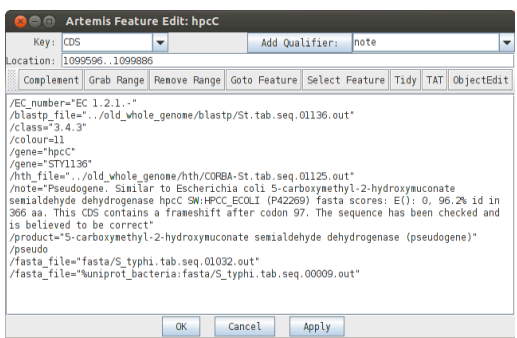
Gene in database

Can you see where the stop codon has been introduced into the sequence of our gene of interest? Search for the highlighted amino acid sequence in *hpcC*. Have a look if you can find the subsequent amino acids of the database hit in any of the three reading frames. You will see the sequence can be found in the second frame! What has happened? The last amino acid in common is a K then the amino acids start to differ till the stop codon. The amino acid K is coded by AAA. The next base is an A, too. This little homopolymeric region can cause trouble during DNA replication if the polymerase slips and introduces an additional 'A'. This shifts the proper reading frame into the second frame.

To correct the annotation we have to edit the CDS now. Left click on the right amino acid continuing the amino acid sequence on the second frame (have a look in the fasta results and look at the sequence of the gene in the database when you are not sure) and drag till the end of the gene. Then click 'Create' 'Feature from base range' and 'OK'. A new blue CDS feature will appear on the appropriate frame line.



As the original gene annotation is too long we have to shorten it. Click on the original *hpcC* CDS, 'Edit' 'Selected features in Editor'. A window will pop up and you can change the end position in 'location' (the end position is the last base of the stop codon).



The new CDS feature can then be merged with the original gene as shown below (1-3).

A small window will appear asking you whether you are sure you want to merge these features. Another window will then ask you if you want to 'delete old features'. If you click 'yes' the CDS features you have just merged will disappear leaving the single merged CDS. If you select 'no' all of the three CDS features (the two CDSs you started with plus the merged feature) will be retained.

2 Click 'Edit'

3 'Selected Features' 'Merge'

1 Select both the original gene-model and the new CDS feature, which is to be merged with it to form a new gene

The screenshot shows the Artemis Entry Edit interface for S_typhi.dna. It displays a genomic map with various features like hpcA, hpcB, hpcC, hpcD, hpcE, hpcF, hpcG, hpcH, hpcI, hpcJ, hpcK, hpcL, hpcM, hpcN, hpcO, hpcP, hpcQ, hpcR, hpcS, hpcT, hpcU, hpcV, hpcW, hpcX, hpcY, hpcZ. A context menu is open over a selected feature, showing options like Duplicate, Merge, Unmerge, etc. The 'Merge' option is highlighted. Below the main window, a 'result' window shows the merged feature details: Selected feature: bases 1461 amino acids 486 hpcC (/EC number="EC 1.2.1.1" /colour=11 /gene="hpcC"/gene="STY1143").

Tip: To select more than one feature (of any type) you must hold the shift key down.

You have now corrected the annotation of the gene. If there is some time left: Is there anything more to correct in this gene? You might need to run another blast search to find out about this.