

Pathogen Genomics: Introduction to genome sequencing and analysis

Adam Reid & Steve Doyle
Wellcome Sanger Institute/LSHTM

LSHTM Pathogen Genomics 2021

Summary

- What is the point of a genome sequence?
- Genome sequencing technologies
- Sequence data files
- Viewing genomes
- Computer practical 1: Viewing genome sequences
- Computer practical 2: Analysis of sequence variants

Why do genome sequencing?

- Reference for molecular biology
 - *Tropheryma whipplei* causes the potentially fatal Whipple's disease. Could not easily be grown. Genome revealed it had lost genes involved in producing amino acids.
- Identify all the genes that determine the function of the organism
 - *Neisseria meningitidis*, a major cause of meningitis. The first vaccine for a particular form of meningitis was identified by looking for candidates in its genome.
 - *Rickettsia prowazekii* is the cause of epidemic typhus, which killed millions in the early 20th century. It cannot reproduce outside of these cells. It was found to have just over 800 genes.
- Examine evolution by comparative genomics
- Track spread of pathogens
- Identify antimicrobial/drug resistance genes and drug targets
 - Mtb researchers made bacteria resistant to a new drug. Genome sequencing identified the gene involved in resistance.
- Basis for other omics technologies – RNA-seq, ChIP-seq, Methylome etc.

Why do genome sequencing? - Video of Wellcome Sanger Institute researchers

Technology overview

- Sanger sequencing produces ~500bp reads
 - Pros: Highly accurate
 - Cons: Expensive, laborious
 - Uses: High quality reference genomes
- Illumina's sequencing-by-synthesis 75-250bp
 - Pros: cheap, lots of reads (e.g. 500 million per run)
 - Cons: short reads
 - Uses: Resequencing, draft genomes, RNA-seq
- Pacific Biosciences Single-Molecule Real Time (SMRT) reads of 5000bp-40000bp
 - Pros: long reads
 - Cons: Fewer reads than Illumina - ~1 million, low accuracy
 - Uses: Reference genomes



Genome sequencing technologies – Interview with Mike Quail

Sequence data

Fasta

>yneN
TTAATTCGCTTCTCATTTCTGCTCATCGCACAGCAGAAGAATTCTCATGAC
TATTATTCGCAATTGCTCACATGGATAAACTAACATACATAAGATAAAACT
CTGCTCACAGCTGAAGAACCTCCGCTCAGTACTGAAGCACCGACTCTTATTCTCTT
CTCCAGCTGTATTAAGCATGACTGATTACGTTTAACTGGTATCCGCTAAATAA
ACATATTGAAATGCATGCGACCACAGTGAAAACAAAATACGCAAAGAGACA
ACTA
>yeR
ACTAACGGCTGCCACCGATAAATTCAAAAAGAGCATACCTAATATTCAAACTAA
GTGGCATCTTCAATAATATAATTAAAGCCCCATGGAGTACCTGAAGGGCTCA
TCGGTAATTCTTCACTTGTAGGAAATGGTACAGAACATTATAATCTTCA
TTATAAATACATGCCATTATTATTTAAACATAGAGGTGCTGGTATTAA
GGGAAGGTGAGATGAAAAGATAGCTGCTATCATTAATTAGTATTATTATGCTG
G
>emrK
AAATCAGGGATTGACCGATGATTATAGTTCAAGTGGCATAATAAGTCTTCTACTA
ATCTCACAGGGTAAGAATTGATTGCAAAAGCCACGGTTAGTCCTGTTGTTTTT
TGACCTCTTAAATTAGGCTCCAACGTTCTGGATAATGTGCAACACATGACTG
GTTTGATGATGAGAATGATGCTCTTCAATTCAATAATTCTGATACTGAGAAA
GAGAGATAATAGTGGACAGATTCAATAAAAAAACTTCAACAGAAAATACT
>evgA
AATAACATTCTACGCCCTGAGGATTAGTAAGAAGACTTATAGTGCACACTGAAACTAT
AAATCTCGGACAATCCCTGATTATTGTCACATTCTGCGACTATTATA
TGGTATACATTGCGAATTATCTTAAAGGAAGCTCAGATTCTTATTTATTGAGAAA
TGAGATGAGCCTTGTCTGTTACTACAGGGAGAAGGGAGATGCTCTTCTGG
GAATAATCTGACGCCAATTATTGATGACCATCTCTTGTCTGCTCAGCAATTCT
>yfdX
TGGCTGTATTACATTAAATCATGATTACATCGATAATAATGACATCTT
GTGGTATAAGAATAGTCTGCGACAGGAAGCATCTTACAATTGTAAGACTAAA
ACTATCTCTGCGATAACTACATGTAAGATAACCCCTTAAAGATCGCTGCT
CTGATTCTCTTACATGCTCACCCAAATAGTGGCGGGTTCTAAACTTGTAAAGA
ATGAGGTAAGTGAACGTTAATTGCCCAGATGGTACAGCAATTCTGCATCT
C

SAM/BAM

Fasto

0HS34_24228;8:1101:1116:7158/2
ATGCAAGTTTTTACTATAAAAATCATATAAGGATAATANNGATAAACATGANNNNATGTGAATGTAATAA
+
BBBBBBFFFFBFFFF<FFFFFFFBBBBFFFFFFFBBBBBFFFFF!!!!<<<FFFFFFF!!!!<<<FFFFFFF
0HS34_24228;8:1101:1116:7481/2
AATATTAAAAAAATGTTGAAAATCAGGAGCATTGGCGCANNNTGCTGCACCTGNNNNNGCACCAAAGCTGAT
+
BBBBBBFFFFFFFFFFFBBBBBFB</BBBBBF&!!<<</FB/B<!!!!/<<BFFFF/</F/
0HS34_24228;8:1101:1116:8637/2
TGTGTTTTTTATTAAATATTTCTAATAATATAATANNAATAAANNNNNAATATATATATAT
+
BB/BBBBFFF</F<FFFFFF<F/FFBBFFBF//FFFFFFBFBBF!!!!<<</F/FF!!!!<</<</F/BBFF//
0HS34_24228;8:1101:1116:52646/2
CCGAAATTATGCCAGAATTACCGAAGATGAGGAGGAAGGCGNNNGAACGACGANNNNNACGACGACACGAG
+
BBBBBBFFFFFFFFFFFBBBBBFFFFFFFFFFFBBBBBFFFFF!!!!<<BFFFFF!!!!<<FFFFFFF
0HS34_24228;8:1101:1116:52943/2
AATATAAAATCATTGTTGAGCTGTTAATACCAAGCGAACCNACATATGATATNNNNNTTACATGCGTATT
+
BBBBBBFBFBFFFF<FFF/FFFFFFFBBBB</FFFFFFF!!!!<<FFFFFFBFFF!!!!<<FFF/FFFFFFF
0HS34_24228;8:1101:1116:65353/2
ACGAAATAAATAAAAAGCTTAAACCAACGACACNNACATATGTTANNNNCATTTAATATTATT
+
//<BB/FFFFFFF/FFFFFFFBFB</BBBBFFFFFB<<<!/<<BFFF&!!!!<<<FBFFFFFBFB
0HS34_24228;8:1101:1117:7618/2
TCTCACTGTGTTGATAAAAATATAATATAAGCTGAGCTGTTACACNNACACAAACACNNNNACACAAAGACATAC
+
B//<E//<E/B/E//E/B/E//<B//<B//<E//E/E//E//<E//<E//<E//<<FFFFFFF/E//

EMBL

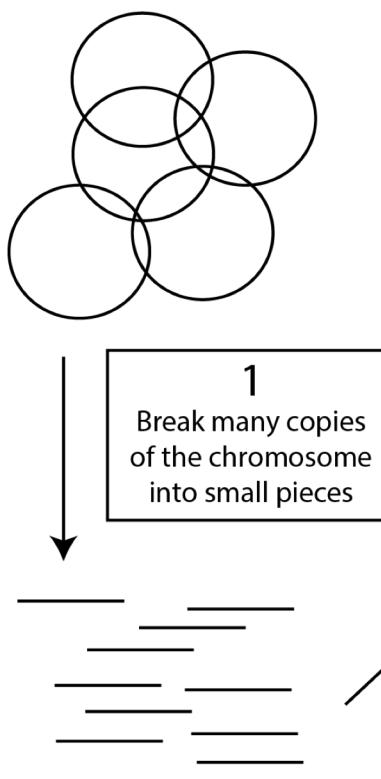
EMBL Flat File

Header

Features (AA seq)

DNA Sequence

What do we do with these data?



3a
Find overlaps between these sequences to make a reference genome

ACTGTCTTCGCG
TACGACTACGAC
GATAGCATACGA
CTTCGGCGCAG
TCTTCGGCGCA
ATAGCATACGAC
CGCAGATAGCAT

2 Sequence the pieces to find out the order of DNA letters in each one

CGCAGATAGCAT
TCTTCGGCGCA ATAGCATACGAC
CTTCGGCGCAG TACGACTACGAC
ACTGTCTTCGCG GATAGCATACGA
ACTGTCTTCGGCGCAGATAGCATACGACTACGAC

3b
Map these sequences to a reference genome for resequencing to find differences

CGCA**T**ATAGCAT
TCTTCGGCGCA ATAGCAT**AG**GAC
CTTCGGCGC**A**T TAGGACTACGAC
ACTGTCTTCGCG **T**ATAGCAT**AG**GA
ACTGTCTTCGGCGCAGATAGCATACGACTACGAC

Annotation (find the genes)

- Look for interesting differences
- Build a tree

How are we going to do our bioinformatics?

- Virtual machine with Linux
- Artemis for viewing genomes
- Various command line tools for mapping, assembling etc.
- Web-based applications

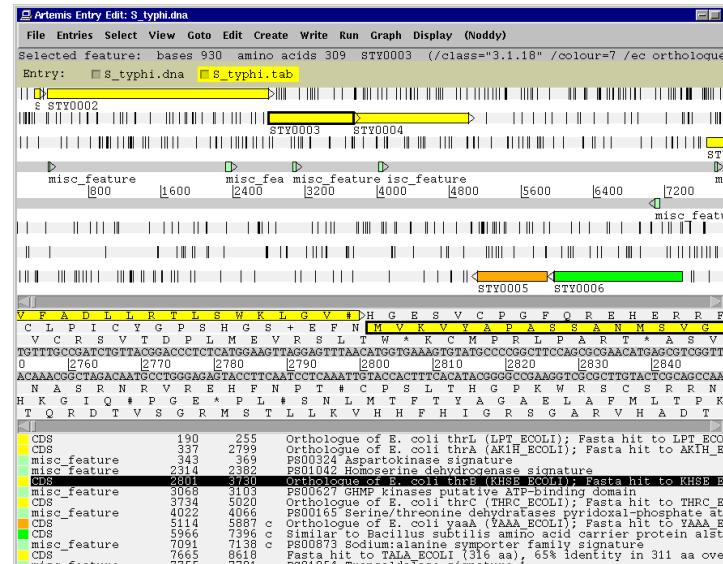
What will we do in the practicals?

- Get familiar with the Virtual Machine
- Computer practical 1: Use Artemis to get familiar with looking at genomes (morning)
- Computer practical 2: Map short-read genome sequencing data to identify differences between closely related bacteria (afternoon)



Genome browser and annotation tool

- visualization of sequence
 - DNA
 - six frame translation
 - Panoramic and sequence view
- Annotation
 - Features
 - Mapped and listed
 - Editable
 - In layers (entry)
- perform and view analysis
 - basic analysis
 - Basic stats & index can be plotted
 - import and view the results of other searches/analysis
 - Different lines of evidence can be seen together





Artemis

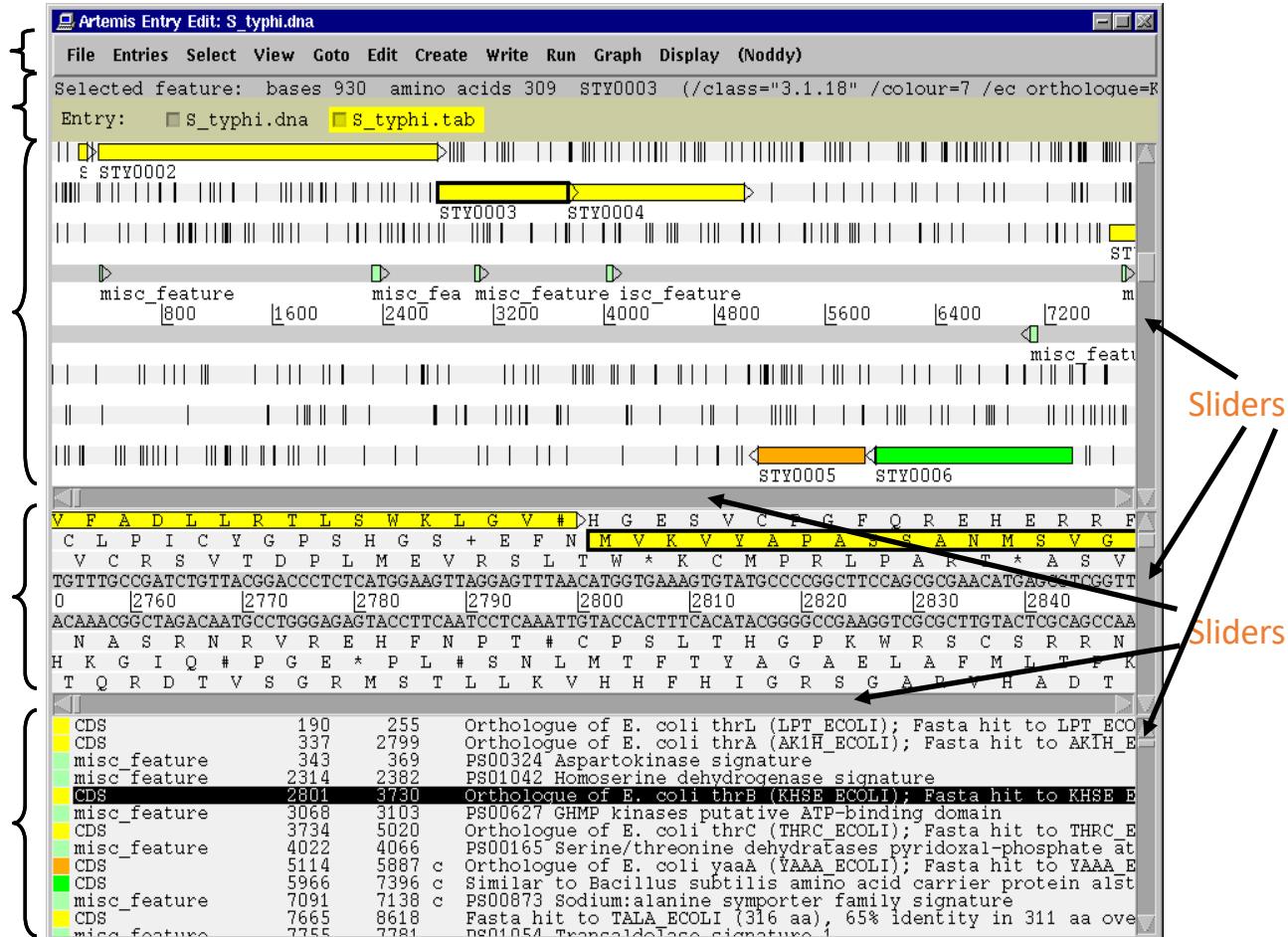
Drop Down Menus

Entry Button Line

Main Sequence View Panel

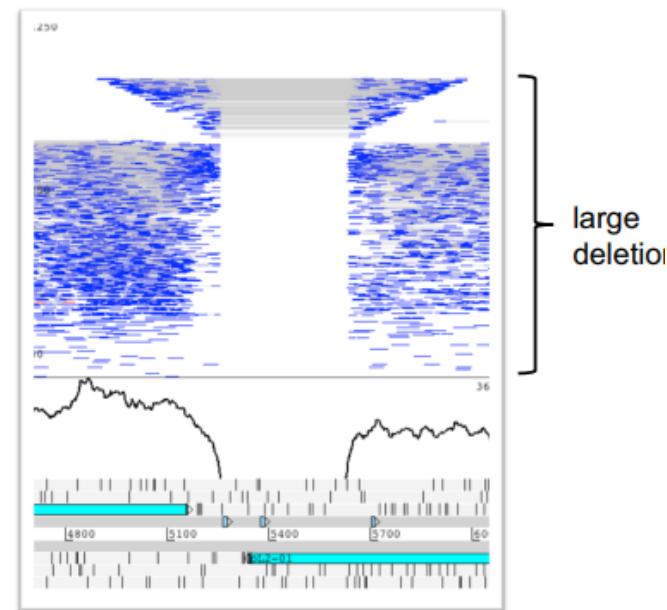
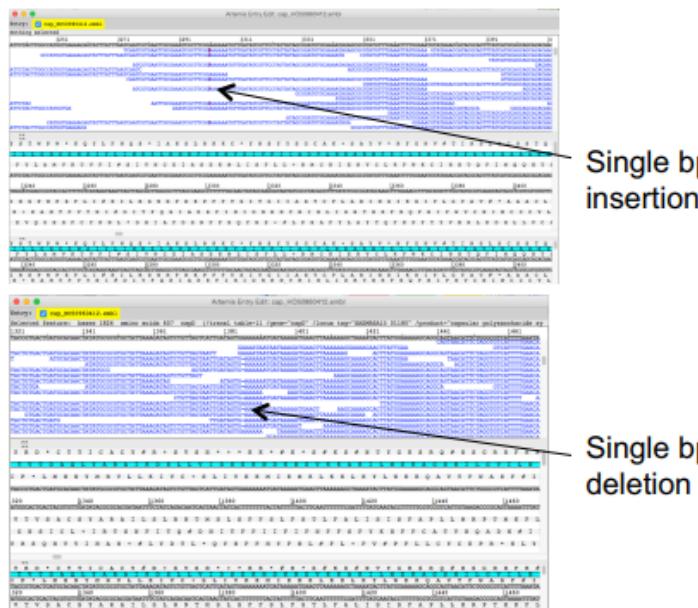
Magnified Sequence View Panel

Feature Menu



Viewing mapped reads in Artemis

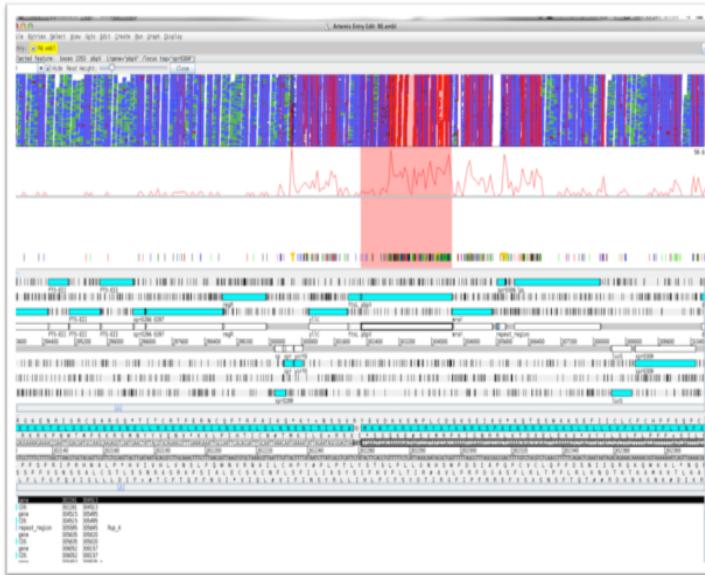
- Illumina data (bam files)
 - identification of indels



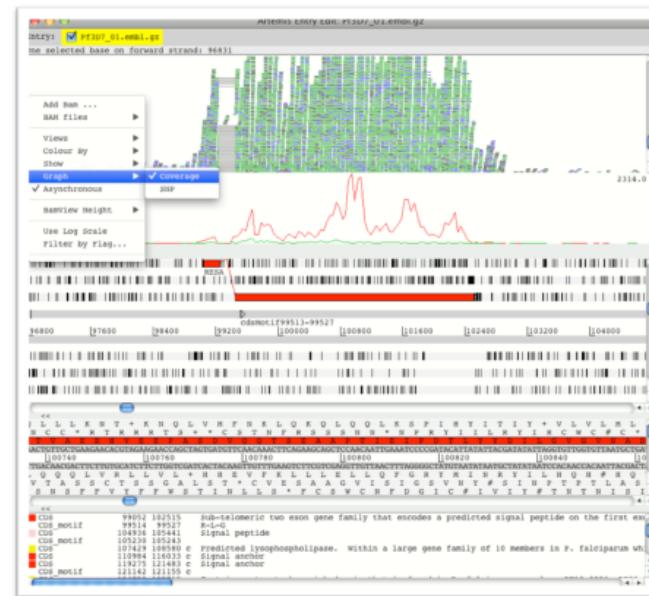
- we will see it on Module 4

Viewing mapped reads in Artemis

Single nucleotide variants (SNVs)



RNAseq data

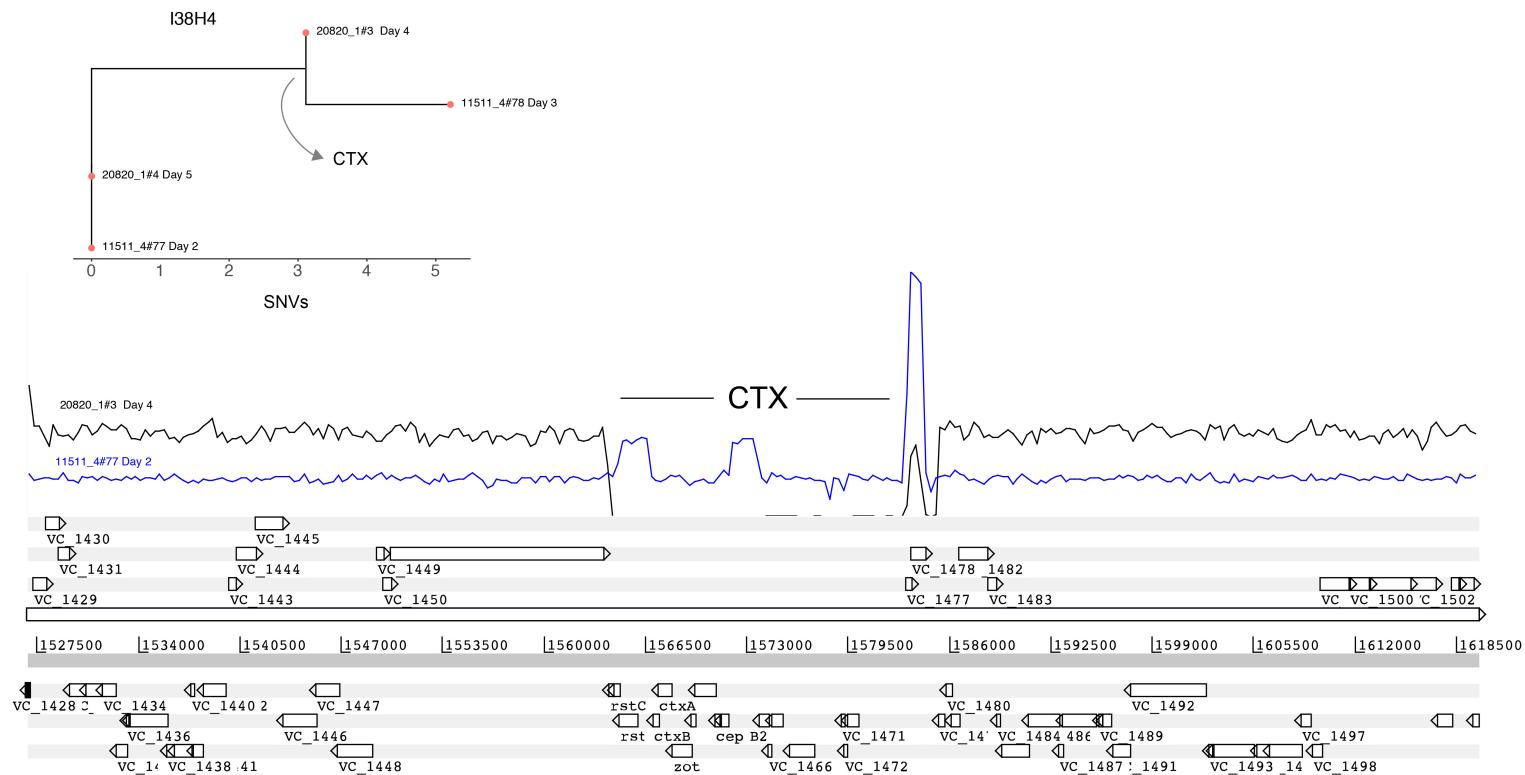


Illumina data (bam files)

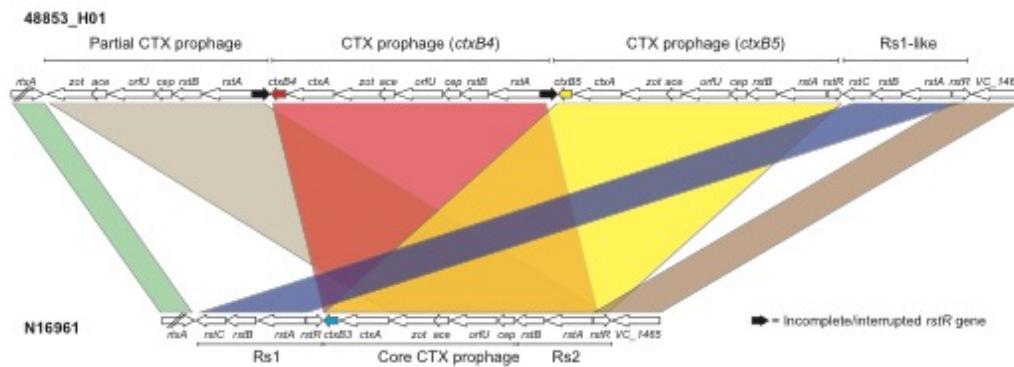
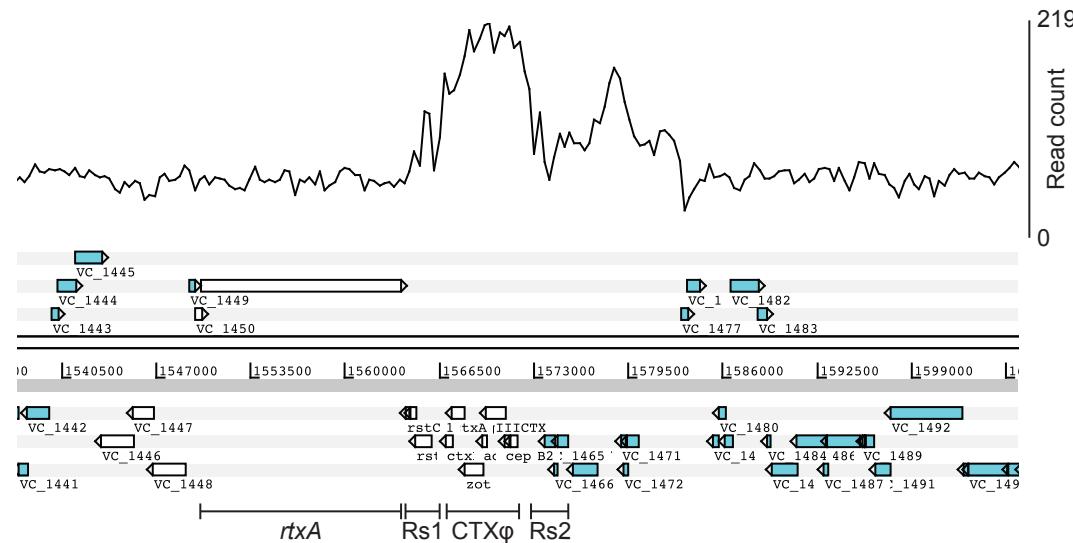
Resequencing and mapping

- Aims to capture information on:
 - Single Nucleotide Variants (SNVs/SNPs),
 - insertions and deletions (indels)
 - Copy Number Variants (CNVs) between individuals of a species.
- As sequences diverge from the reference, mapping becomes progressively less effective

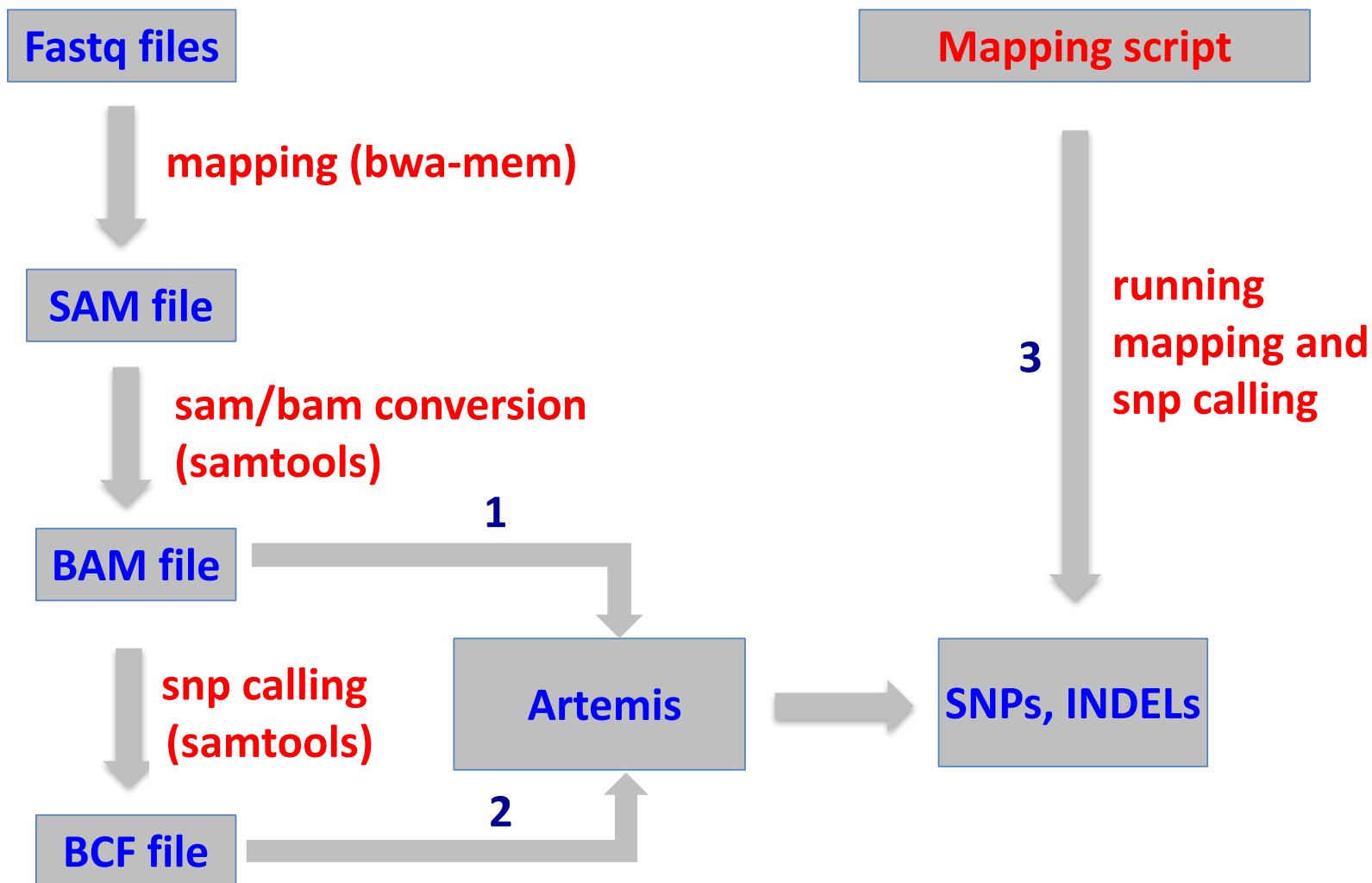
Gene presence / absence



Copy number variation



Mapping sequencing reads-Workflow



The Swedish Story



- Prior to 2006, *C. trachomatis* in Sweden was following the same pattern as in the UK
- In 2006, across Sweden there was a reported drop in cases

The Swedish story

- It was noticed that counties using the NAATs (Abbott / Roche) diagnostic system showed a drop in *C. trachomatis* cases in 2006
- Counties using other diagnostic methods (BecktonDickinson) still showed an increase in cases, in line with that of previous years
- The obvious conclusion was that the NAATs was missing a subset of infections
- Why?
 - Let's find out using the awesome power of genomic sequencing!!!!

Summary

- Computer practical 1
 - Use Artemis to view genomes
 - Understand genome data files
 - Understand relationship between the sequence and annotation
 - Understand what bacterial genomes look like and how they are arranged
- Computer practical 2
 - Practice short read alignment
 - View mapped reads in Artemis
 - Call and view SNPs
 - Uncover why the PCR test failed for new variant Chlamydia