

Pathogen Genomics: Introduction to genome sequencing and analysis

Adam Reid & Steve Doyle
Wellcome Sanger Institute/LSHTM

LSHTM Pathogen Genomics

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 for 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

TTAATGCGCTCTTCTCATTTCTTCTGCTGTCATCCGCACAGCAGAAGAATTCTCATTTGAC
 TATTATTTTCGAATTTGCTCACATGGATTAAATTAACTACATACATAAGATATAAACT
 TCTGCCTACAGCTGTAGAAGAACTCCGCTCAGTACTGAAGCACCAGTCTTATTTCTCTTT
 TCTCCAGCCTGTATATTAAGCACTGATTAAACGATTTTAAACGTTATCCGCTAAATAA
 ACATATTTGAAATGCATGCGACCACAGTGAAAAACAAAATCACGCAAAGAGACAACATA
 A
 >yegR
 ACTAACGGCTGCCACCATAATTTCAAAAAAGAGCATATACCTAATATTCAACTAAACA
 GTGGCATCTTCAATATAATATATTAAAGCCCCCATGGAGTTACCTGAAGGGCCTCAATG
 TCCGTAATTCCTACTTATGTAGGAAATGTTGTACAGAACATTTATTATAATCTTATCTT
 TTATAAATCATGCCATTATATATTAAACACTAGAGAGTGTCTGTTGGTATTTAATGG
 GGGAAAGTGAGATGAAAAAGATAGCTGCTATATCATTAATTAGTATTTTATTATGTCTG
 G
 >emrK
 AAATCAGGGATTGTACCGGATGATTTATAGTTTCAAGTTGGCATAAGCTTCTTACTA
 ATCTCACAGGCGTAAGAATTGTATTGCAAAAGCCACGGTTTAGTCCTCTGTTGTTTTTT
 TCGACCTCATTTAAATTAAGCTCCCAACGTTCTCGGATAATGTGCACAACATGCACGTGT
 GTTTGATATGAAGATGAATGCTCTTTTCATTCAGTTTCTATAAATTTTCATCTAGAAAT
 GAGAGATAATAGTGAACAGATTAATTCAAATAAAAAACATTTCAACAGAAGAAAATACT
 T
 >evgA
 AATACAATTCTTACGCCTGTAGGATTAGTTAGAAGACTTATAGTGCCAACTTGAACATAT
 AAATCATCCGTACAATCCCTGATTTTATTGTTGACATTTCTATTATGCCGACTATTTATA
 TGGTATACTTGTGCAATTCTTTAAAGGAAGCTCAGATTTCTTATTTTATTGAGAAA
 TGAGATGACGCCATTATGCTGTATTACTACAGGAGAGGGAATGTTCTATTGCAAGG
 GAATAATCTATGAACGCAATAATTATTGATGACCATCCTCTTGCTATCGCAGCAATTCTG
 T
 >yfdX
 TGCGTGATTTACATTTAATTAATCAGTATTTACATCGATATAATAAATGACATCTCTTT
 GTGGTATATAAGAAATAGTTCTCGCAGCAGGAACATTTCTACAATTGTAGACATAA
 ATACTCTTTCGCAATTAACCACTGTAAGATAACCTTTCAAAATGACCGTGTGCTCT
 CTGATTTCTCATTTTCATGCTCACCAATATGATGGCGGGCTTTCTAAAACATGTTAAAGA
 TAGGGTAAGTATGAAACGTTAATTATGGCCACGATGGTCAACAGCAATTCCTGGCATCTT
 C

SAM/BAM

[illegible]

Fastq

#H\$34_24228;8;1101;1116;7158/2
ATGCAGATTTTTTACTATAAAAARTCATCATAAGGATAANNNGATAACAATGANNNNNNATGTGAATCTAATAA
+
BBBBBBFFBFFFF<FFFFFFFFFBFFFFFFFFFFFFFFF!!<<<FFFFFFF!!!!<<<FFFFFFFFF
@H\$34_24228;8;1101;1116;7461/2
AATATTAAAAAAATGGTTGAAATCCGACAGCATTTCGCCANNTGCTGCACCTGNNNNNGACCCAAGCTGAT
+
BBBBBBFFBFFFFFFFFFBFFFFFFFFFB/FBBBBBF<F!<<</<F/B/<!!!!//<<<BFFFF/<F/
@H\$34_24228;8;1101;1116;8637/2
TGGTGTTTTATTATTATTAAATATATTCTAATAATAATATANNNATAAATAAAAAANNNNNNAATATATATATAT
+
BB/BBFFFF<F<FFFFFFF<F/FBBBBBF/FFFFFFBFBF!!<<<<<F/FF!!!!<</</<F/BFFF/<
@H\$34_24228;8;1101;1116;52646/2
CCGAATTATTGCCTAGAATTCACGAAGTAGGAGGAGGCGNNGGAACACGACGANNNNNACGACCACCAAGC
+
BBBBBBFFBFFFFFFFFFBFFFFFFFFFBFFFFFFFFF!!<<BFFFFFFF!!!!<<FFFFFFFFF
@H\$34_24228;8;1101;1116;52943/2
AATATAAACATTGATTGAAGTGTATTAAACACAGCGACNNACATATGATATNNNNNTTACATCGCTATT
+
BBBBBBFBFBFFFF<FF/FFFFFFFFFFFFF/<FFFFFFF!!<<FFFFFBFFF!!!!<<FF/FFFFFFF
@H\$34_24228;8;1101;1116;65353/2
ACGAAATAAATAAAAGGTATTTAAACCAAAAATGATAATANNCAATATGTTTANNNNNCATTTAATATTATT
+
//<BB/FFFFFFF/FFFFFFFBFBFB<BFFFFFFFB<<<!!!!<<<BFFB//!!!!<<<FFBFBFBFBF
@H\$34_24228;8;1101;1117;7618/2
TGCTCACCTGTTGATAAAAAATATAATAAGTTAGAGTTACANNACACACACACNNNNNCAACAGAGGCATAC
+
B//<F/////B/B/FF/FF/F/B//BFF///<B//B//!!!!<<F/FF/F/!!!!<<<FFFFFFFFF/F

EMBL

```

AC AF115338 standard; DNA; PRO; 591 BP.
AC AC115338;
SV AF115336.1;
DT 03-JUN-1999 (Rel. 59, Created)
DT 23-AUG-1999 (Rel. 60, Last updated, Version 2)
DE Pseudomonas fluorescens ECF sigma factor SigM (sigK) gene, complete cds.
HG
OC Pseudomonas fluorescens
OC Bacteria: Proteobacteria; gamma subdivision: Pseudomonadaceae: Pseudomonas.
RN 1-591
RX MEDLINE: 95669842.
RI Brinman F.-J., Schoofs G., Hancock R.E., De Mot R.;
RT "Influence of a putative ECF sigma factor on expression of the major outer membrane protein, OmpF, in Pseudomonas aeruginosa and Pseudomonas fluorescens";
RL J. Bacteriol. 181(16):4746-4754(1999).
RN [2]
RF 1-591
RA De Mot R.;
RJ ;
RL Submitted (04-DEC-1998) to the EMBL/GenBank/DDBJ databases.
AB Pseudomonas fluorescens; Genet. res. Applied Plant Sciences, SPRETEL Merckelstraat 52, Heverlee B-3001, Belgium
DR SPTREML: QXKAL7; QXKAL7.
FH
Key Location/Qualifiers
FT source 1..591 /db_xref=taxon:294
FT /organism=Pseudomonas fluorescens
FT /strain=K14.1
CDS 1..591 /codon_start=1
FT /db_xref=SPTREML:QXKAL7
FT /transl_table=1
FT /gene=sigM
FT /product="ECF sigma factor SigM"
FT /protein_id=AAD34329.1
FT translation=MGGGGLTGVYVDELSEELVARSHTSLTYNVEITVYCEIQRYSQ TLFNQCYAYLRNGDGDQVWGLRQLGKSGEQLPQSTFLYSTIVYNTCTQPKS
FT RRKRRLMDLLPDLA SEAGL PQEWGLDGGLLVYVPNI DQGLVLPVAFLEPFE
FT IADLMHSLATVAVLQGLDGLDTPFASVTET
SQ
Sequence 591 BP; 151 A; 133 C; 170 G; 131 T; 0 other;
```

EMBL Flat File

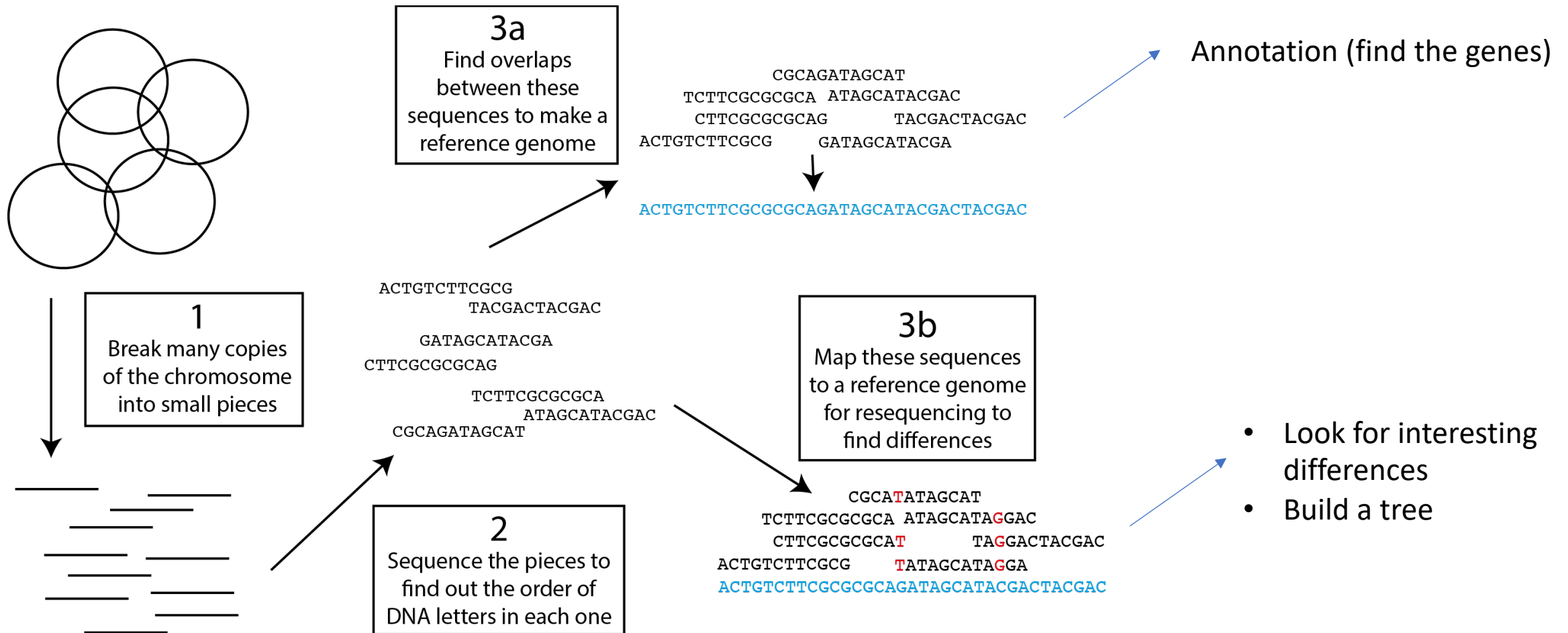
Header

- Title
- Taxonomy
- Citation

Features (AA seq)

DNA Sequence

What do we do with these data?



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

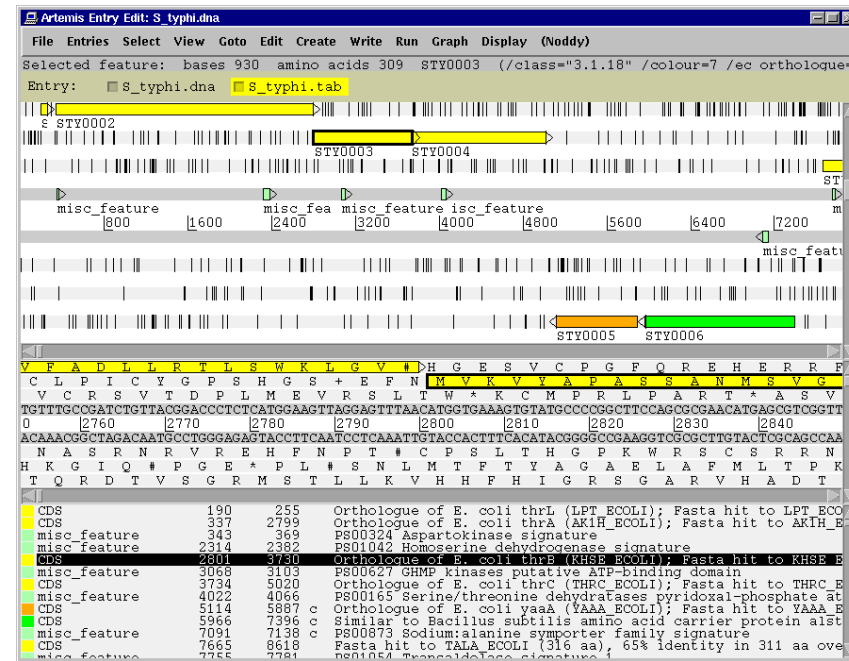
What will we do in the practicals?

- 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

Artemis Entry Edit: S_typhi.dna

File Entries Select View Goto Edit Create Write Run Graph Display (Noddy)

Selected feature: bases 930 amino acids 309 STY0003 (/class="3.1.18" /colour=7 /ec orthologue=K)

Entry: ☐ S_typhi.dna ☒ S_typhi.tab

STY0002 STY0003 STY0004

misc_feature misc_feature misc_feature misc_feature

800 1600 2400 3200 4000 4800 5600 6400 7200

STY0005 STY0006

V F A D L L R T L S W K L G V # H G E S V C P G F Q R E H E R R F

C L P I C Y G P S H G S + E F N M V K V Y A P A S S A N M S V G

V C R S V T D P L M E V R S L T W * K C M P R L P A R T * A S V

TGTTGCCGATCTGTTACGGACCCTCTCATGGAAGTTAGGAGTTTAAATGTTGAAAAGTGTATGCCCCGGCTTCCAGCGCGAACATGAGCGTTCGGTT

0 2760 2770 2780 2790 2800 2810 2820 2830 2840

ACAAACGGCTAGACAATGCCTGGGAGAGTACCTTCAATCCTCAAATGTACCACTTTACATACGGGGCCGAAGGTCGCGCTTGTACTCGCAGCCAA

N A S R N R V R E H F N P T # C P S L T H G P K W R S C S R R N

H K G I Q # P G E * P L # S N L M T F T Y A G A E L A F M L T P K

T Q R D T V S G R M S T L L K V H H F H I G R S G A P V H A D T

CDS	190	255	Orthologue of E. coli thrL (LPT_ECOLI); Fasta hit to LPT_ECOLI
CDS	337	2799	Orthologue of E. coli thrA (AK1H_ECOLI); Fasta hit to AK1H_ECOLI
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
misc_feature	3068	3103	PS00627 GHMP kinases putative ATP-binding domain
CDS	3734	5020	Orthologue of E. coli thrC (THRC_ECOLI); Fasta hit to THRC_ECOLI
misc_feature	4022	4066	PS00165 Serine/threonine dehydratases pyridoxal-phosphate at
CDS	5114	5887	Orthologue of E. coli yaaA (YAAA_ECOLI); Fasta hit to YAAA_ECOLI
CDS	5966	7396	Similar to Bacillus subtilis amino acid carrier protein alst
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 ove
misc_feature	7755	7781	PS01054 Transaldolase signature 1

Sliders

Sliders

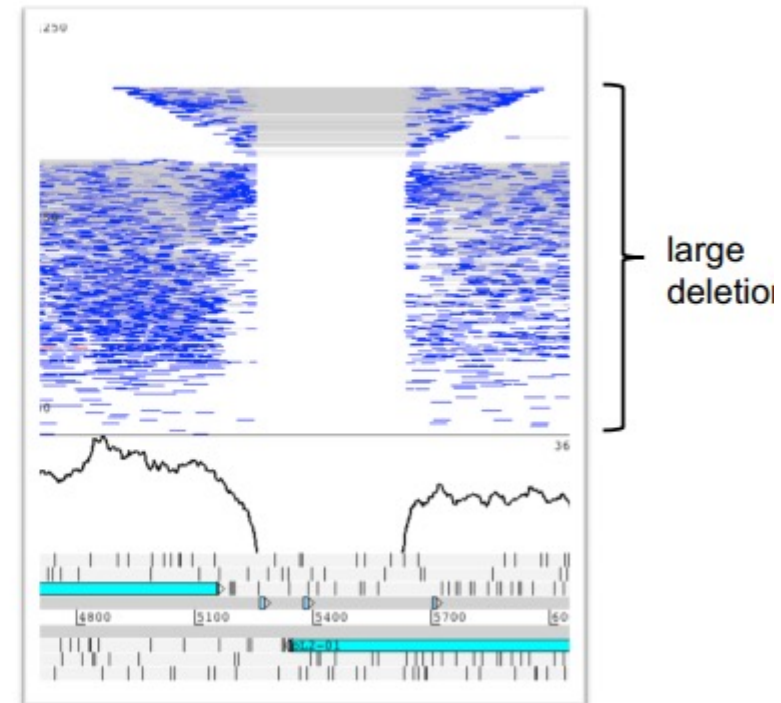
Viewing mapped reads in Artemis

- Illumina data (bam files)
- identification of indels



Single bp
insertion

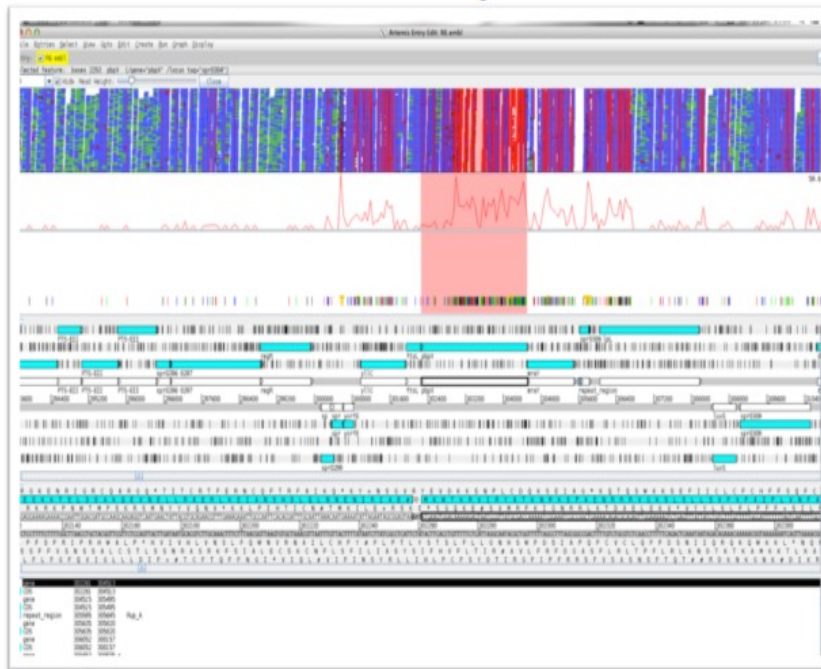
Single bp
deletion



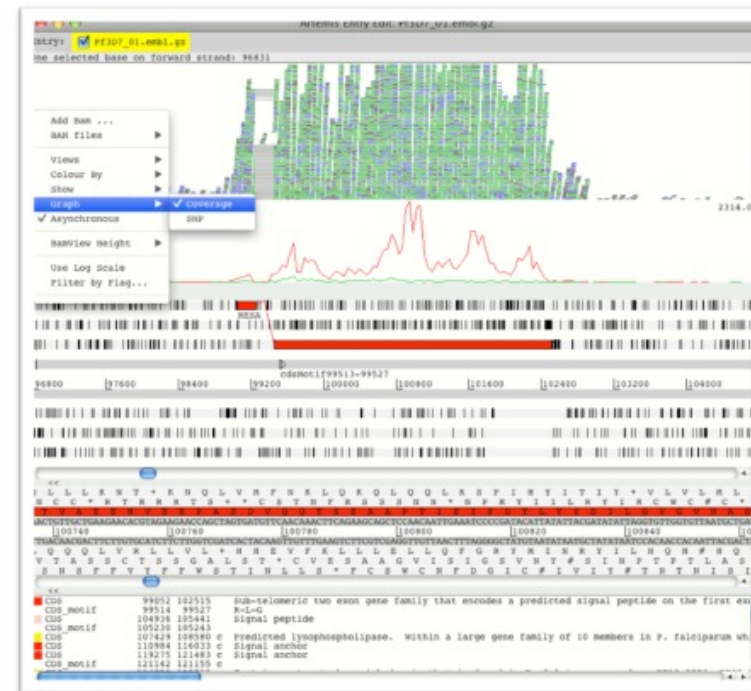
- we will see it on Module 4

Viewing mapped reads in Artemis

Single nucleotide variants (SNVs)



RNAseq data

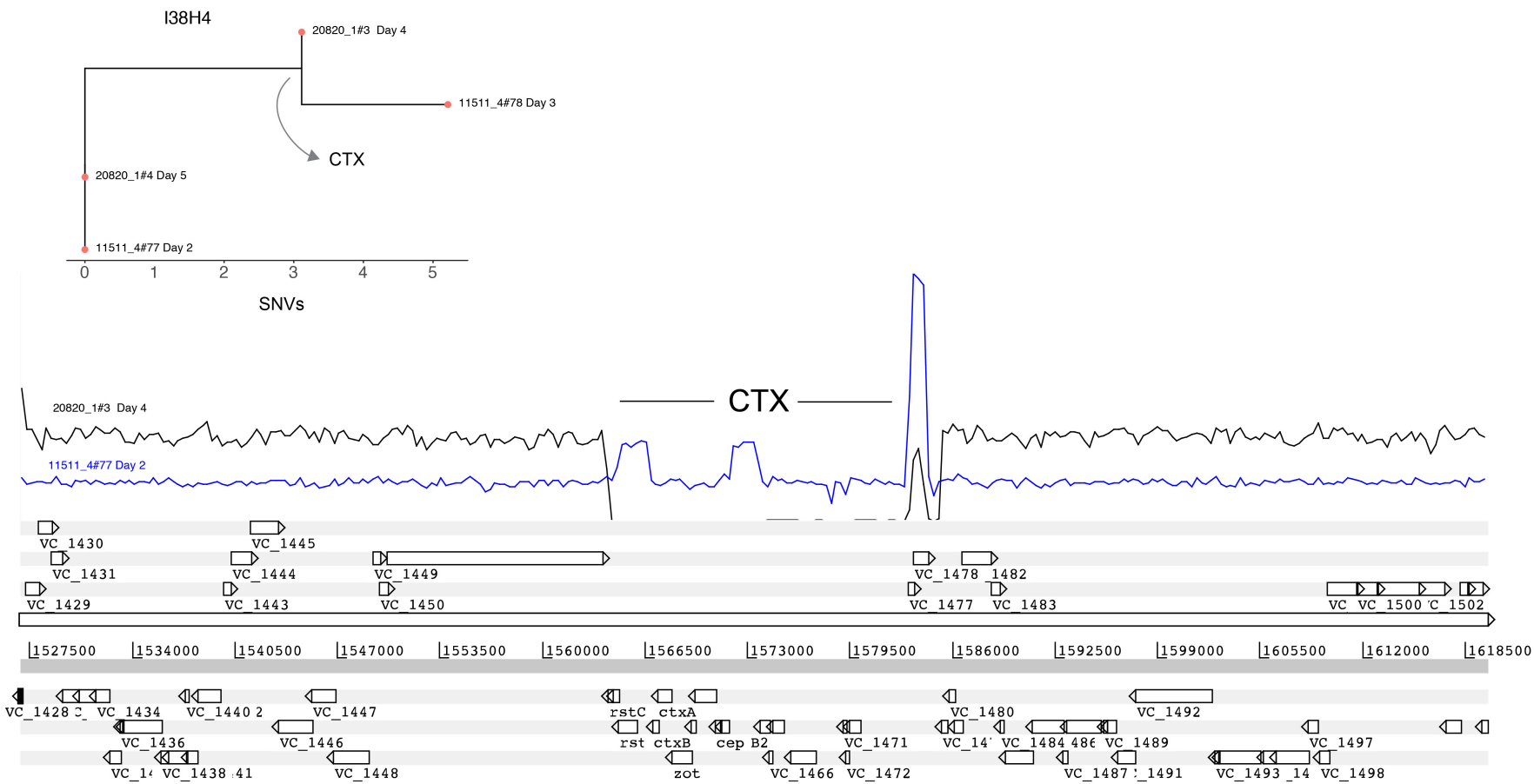


Illumina data (bam files)

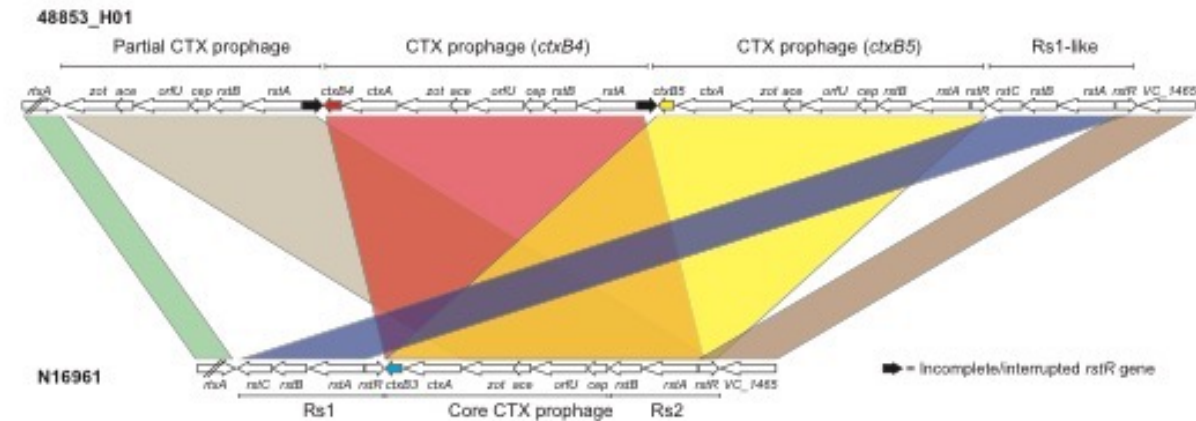
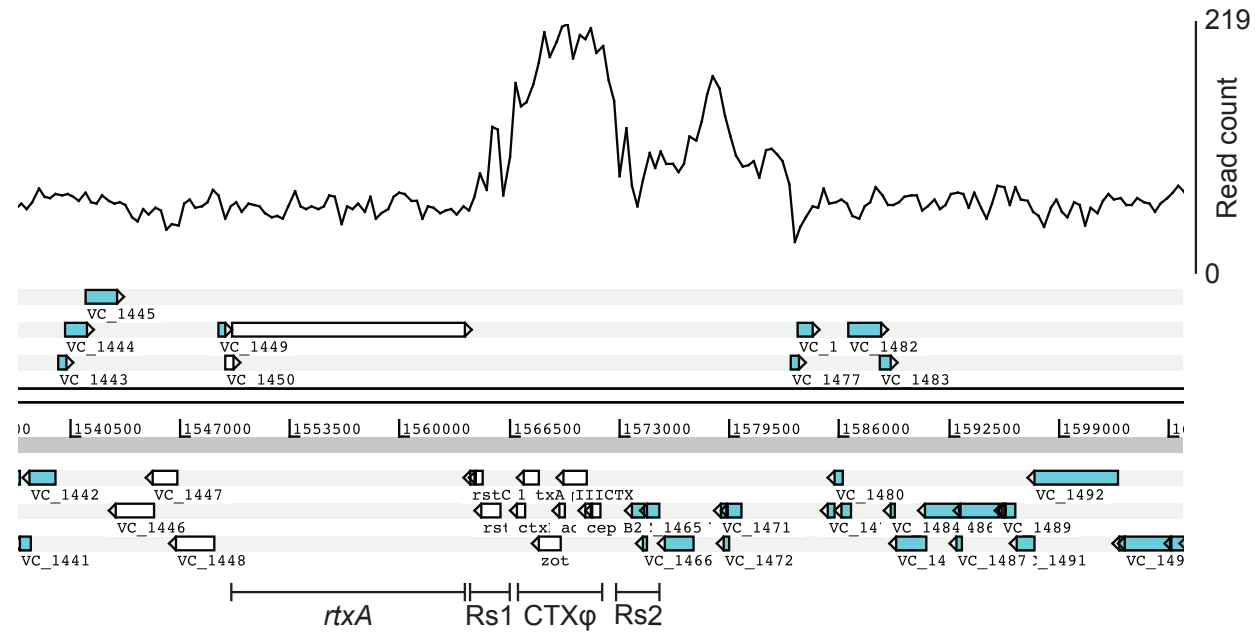
Resequencing and mapping

- 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

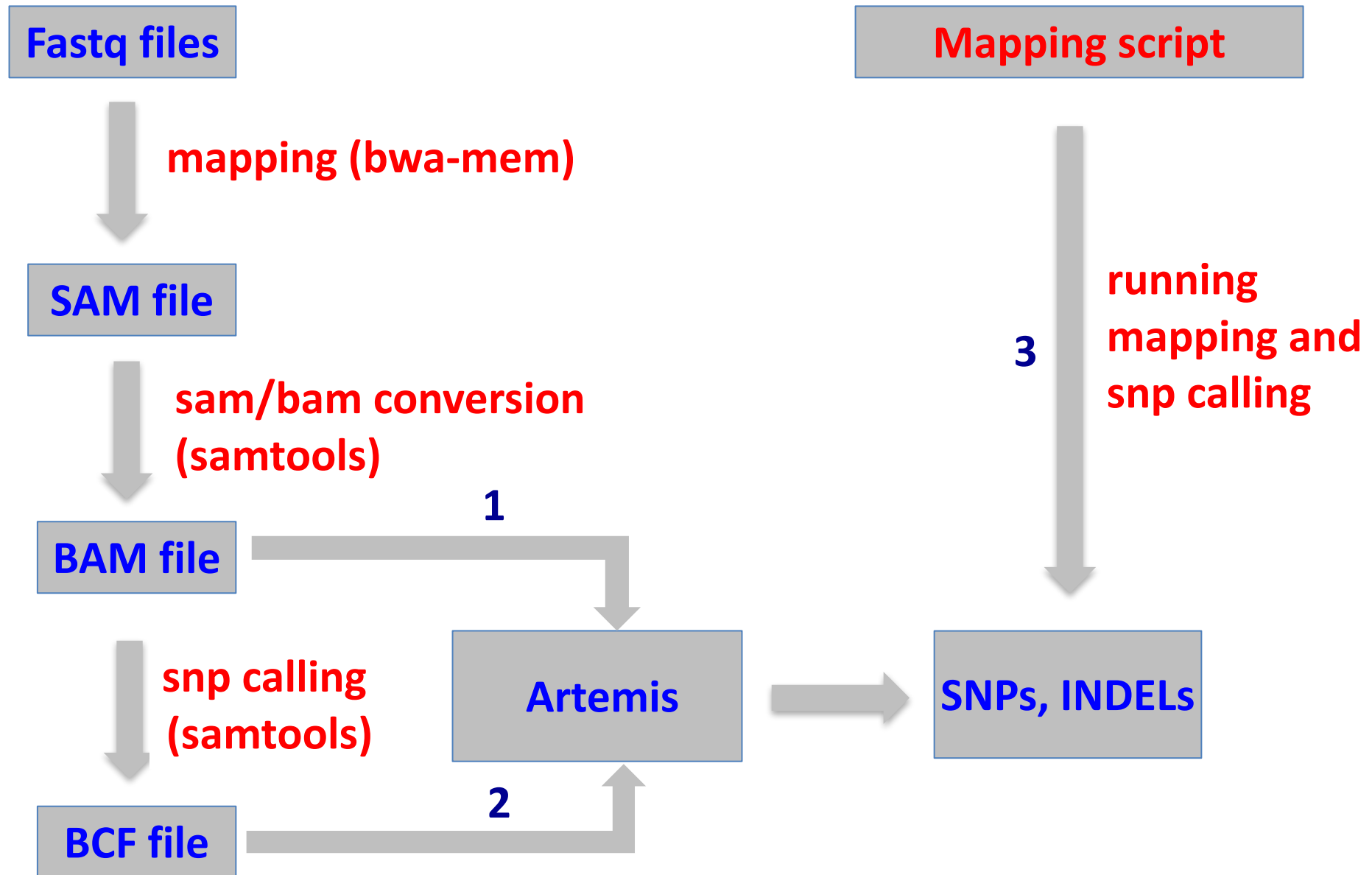
Gene presence / absence



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