

Pathogen Genomics – working with pathogen genomes

Adam Reid

Wellcome Sanger Institute

LSHTM Pathogen Genomics 2020

Summary

- What is the point of a genome sequence?
- Genome sequencing technologies
- Sequence data files
- Viewing genomes

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

tyneN
TATTAGGCTCTCTCATTTCTTCTGCTGTCATCCGCACAGCAGAAGAATTCCTCAITGAC
TATTATTTCGCAATTTGCTCACATGGATTAATTAACACTACATACTATAAGATATAAACT
TCTCGCCTACAGCTGTGAAGAACTCCGCTCAGTACTGAAGCACCAGCTCATTTCTCTT
TCTCCAGCTGTATATTAAGCATACITGATTAACGATTTTAAAGCTTATCCGCTAAATCA
ACATATTTGAAATGCATGCGACCACAGTGAAAAACAAATCACGCCAAAGAGACAACATA
A
>yegR
ACTAACGGCTGCCACCATAAATTTCAAAAAAGAGCATATACCTAATATTTCAACTAAACA
TGGCATCTCTCAATATAATATAAAGGCCCATCGGATGACCTGAAGGGCCCTCAATG
TCCGTAATTCCTACTTATGTAGGAATGTTGTACAGAACATTTATATAATCCTATTCAA
TTATAATAATCATGCCATTATTATATTAAACACTAGAGAGTGTCGTGGTATTTAATGG
GGAAGGTGAGATGAAAAGATAGCTGCTATCATTAATGATGATTTTATTATGTCTG
G
>emrK
AAATCAGGGATGTCACCGATGATTATAGTTTCAAGTTGGCCATATAAGTCTTTTACTA
ATCTCACGGCGTAAGAATGTATTGCAAAAGCCAGCGTTAGTCTCTGTTGTTTTTT
TGCACATCTTAATTAATAGGCCCTCAACGTCTCTGGGATAATGTGCAACACAGCACTGT
TTTGTATGAAGAATGAATGCTCTTTTCTTCAATCAATATAATTTATCTATGAGAA
GAGAGATAATGTGGAACGATTAATTCAAATAAAAAACTCTCAACGAAGAAAATCT
T
>evgA
AATACAATCTTACGCCTGTAGGATTAGTAAGAAGACTTATAGTGCCAACTTGAAACTAT
AATCATCGGTACAATCCCTGATTTTATGTTGACATTTTATTATGCCGACTATTIATA
TGGTATACTTGTGCAATTATCTTAAAGAGACCGTCAGATTTTCTATTTTATTGAGAAA
TGAGATACGCCCTTATGTCTGATTACTACAGGGAGAAGGAGAGGCTCTTATGCAAAAGG
GAATAATCTATGAACGCAATAATTATTGATGACCATCCTCTTGCTATCGCAGCAATTCGT
>yfdX
TGGCTGTAATTACATTTAATAATCATGATTTACATCGATATAAATGACATCTCTTT
GTGGTATAAGAAGTAGTTCTCTCGCAGGAAGACATTTCTCAATGTGAAGACATAA
ATACCTCTTCGCAATAATCAACAATGTGAAGATAACCTTTCAAAATGACCGTGTCTCT
CTGATTTCTCATTTCACTGCTCACCACATGATGGCGGCGTTTCTAAACTGTTAAAGA
ATAGGTAAGTATGAAACGTTTAAATTATGGCCACGATGGTACAGCAATTCGCGCATCTT
C

SAM/BAM

[illegible]

Fastq

```

@HS34_24228:8:1101:1116:7158/2
ATCGAGTATTTTTTACTATAAAAAATCATCATAGGATAATANNCATACAACTAGNNNNNATGTGAATGTATAAA
BBBBBFFFFBFFFFFFF<FFFFFFFBBFFFFFFF<FFFFFFF<!!!!<FFFFFFF<!!!!<FFFFFFF<!!!!
@HS34_24228:8:1101:1116:7461/2
AATATTAAAAAATGTTGAAATCCACAGGCAATGCCGCANNTGCTGCACCTGNNNNNNGCAGCCAAAGCTGAT
BBBBBFFFFFFFFFFFFFFFFFFFFFFFFF<B/FBBBBFB<!!!!<<B/FB<!!!!<<B/FFFF<<F/
@HS34_24228:8:1101:1116:8637/2
TGCTGTTTTTATTATTTTAAATATATCTAATAATAATANNATAATAAAAAAANNNNNNAATATATATATATAT
BB/BBFFFF<F/FFFFFFF<F/FFBFFBF/FFFFFFBFBF!!!!<<F/FF!!!!<<F/BFFFF<
@HS34_24228:8:1101:1116:52646/2
CCGAATTATGCCTAGAATTCACAGATGAGGAGCAGCGCNGGAACACAGNANNNNNACAGACCACGAG
BBBBBFFFFFFFFFFFFFFFFFFFFFFFFF<BFFFFFFF<!!!!<FFFFFFF<!!!!
@HS34_24228:8:1101:1116:52943/2
AATATAAAACATGATTGAGGTGATTAACCAAGCGAACNNACATATGATATNNNNNTACATCGGTATT
BBBBBFBFBFFFF<FF/FFFFFFF<FFFFFFF<FFFFFFF<!!!!<FFFFFFF<!!!!<FF/FFFFFFF<
@HS34_24228:8:1101:1116:65353/2
ACCAATAATAAAAAAGTATTTAAACCAAAAATGATAATANNCAATATGTTTANNNNNACATTATATTATT
//<BB/FFFFFFF/FFFFFFFBFBFB<BFFFFFFFB<!!!!<<BFBF<!!!!<<FFBFBFBFB
@HS34_24228:8:1101:1117:7618/2
TGCTCACTGTGTGATAAAAAATATATAAGTTAGAGTTACANNACACACACNNNNNACACAGAGCAATAC
B/<F<F/!!!!<B/BF/FF/F/B/B/FF/!!!!<B/<B/<!!!!<F/FF/<F/!!!!<F/FFFFFFF/F/

```

EMBL

```
ID AF015338 standard: DNA; FEO: 591 BP.
AC AF015339;
SV AF015339.1;
DT 03-JUN-1999 (Rel. 9, Created)
ZS 23-JUN-1999 (Rel. 60, Last updated, Version 2)
NC Pseudomonas fluorescens ECM surface factor SigaK [sigK] gene, complete cds.
MG
CC Pseudomonas fluorescens
CO Bacteria; Proteobacteria; gamma subdivision; Pseudomonadaceae; Pseudomonas.
RH
PI 1-591
OR NCBI
RS MCLINKED: 9969694.
RA Brittain F.S., Schorffle G.W., Hancock R.E. De Mot R.;
RT "Influence of a Putative ECF sigma factor on expression of the major outer
membrane protein, OprM, in Pseudomonas aeruginosa and Pseudomonas
SM S. J. Bacterial. 181(14):4746-4754(1999).
RL
RM 2
RN 1-591
RA De Mot R.;
RT .
ID Submitted (04-OCT-1998) to the EMBL/Genbank/DDBJ databases.
AU K.A. Jaussens Laboratory of Genetics, Belgium Plant Sciences, N.V.
MC Merckxian ZN, Heverlee B=3001, Reigius
DR SFPRDBML: OXK417 OXK417
CV accession/Qualifiers
FH
FT source 1..591
FT /db_xref=taxon:1214?
FT /organism=Pseudomonas fluorescens"
FT /strain="ML-14"
FT CDS 1..591
FT /codon_start=1
FT /translation="SFPRDBML:OXK417"
FT /transl_table=1
FT /gene="sigK"
FT /product="ECF sigma factor SigK"
FT /protein_id="AB043433.1"
FT /translation="MRGQVSTGLTSLYDPEELSEALVSRHLELVHTVTRAVELHVRQRV
TLNVCAATLITLGRCDLGGVLTGVLGVLYLNGDLSPNTFLIVISITINTEITCYRGE
RHHHMLLNLNLAAGLAPGASLAFETALQVTPPLPQLPQLPVYVMTGTHLPTLSTVFIEE
LITLDHLSLLHFGVLRVQVALLDLKFAETATE"
SQ
Sequence 591 BP: 187 A+ 133 C+ 170 T+ 6 G+ other :
gtatgaataag atcccaacg atcttcacg taagcccgc cggagctttg tatgatagg 60
tattgtgtg gtcgtaccg cgacttttg caaaccttagc agatagaag agatagaag 120
tattgtgtg agatagaag agatagaag agatagaag agatagaag agatagaag 180
cgagactttg tccgtaccg agactgttt ttgactgtt ttgactgtt ttgactgtt 240
agatagaag agatagaag agatagaag agatagaag agatagaag agatagaag 300
tatgatagg taccgtaccg gctgtgttt ttgactgtt ttgactgtt ttgactgtt 360
```

EMBL Flat File

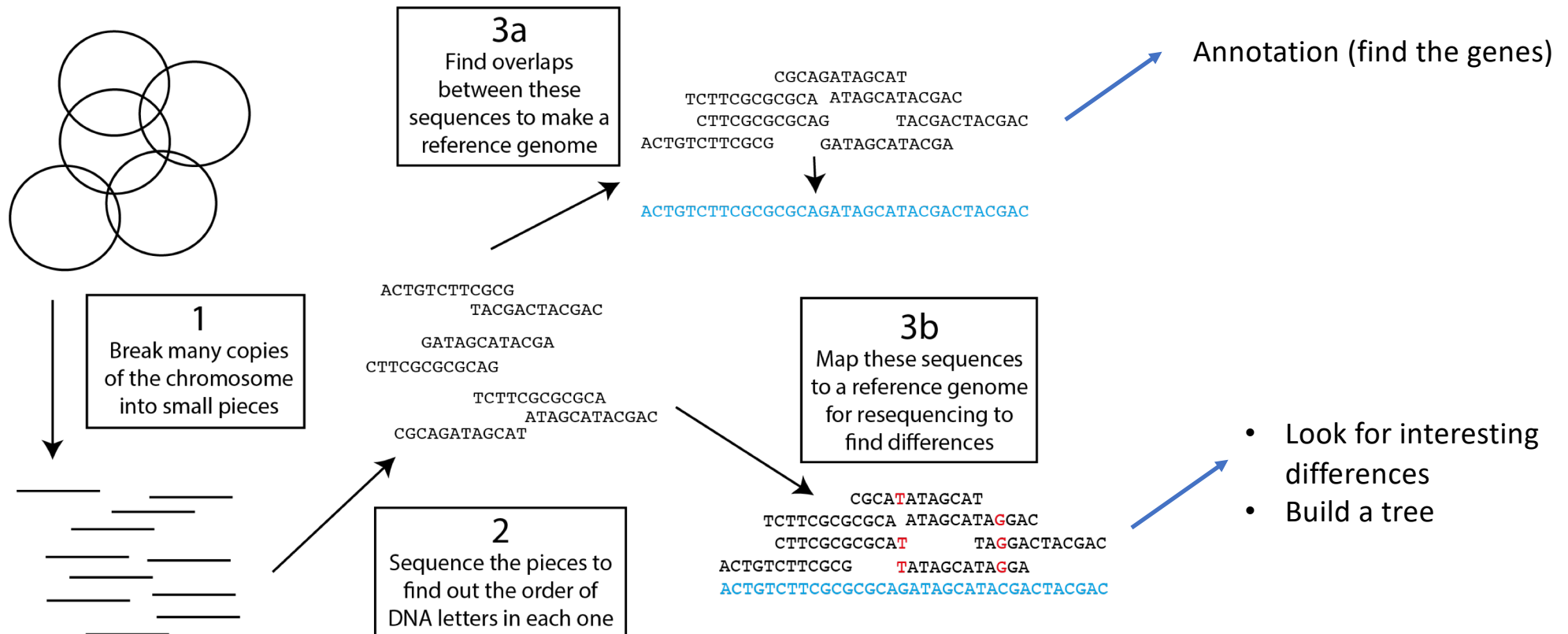
Header

- Title
- Taxonomy
- Citation

Features (AA seq)

DNA Sequence

What do we do with these data?



Doing bioinformatics

- Has anyone here done any bioinformatics?
- 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 today?

- Get familiar with the Virtual Machine
- Use Artemis to get intimate with genomes (morning)
- Map Illumina genome sequencing data to understand differences between closely related bacteria (afternoon)