What can we do with microbial WGS data?

A/Prof Torsten Seemann

Victorian Life Sciences Computation Initiative (VLSCI)
Microbiological Diagnostic Unit Public Health Laboratory (MDU PHL)
Doherty Applied Microbial Genomics (DAMG)

The University of Melbourne

About me

Melbourne, Australia



Microbial genomics + bioinformatics











Microbiological Diagnostic Unit

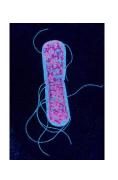
- :: Oldest public health lab in Australia
 - : established 1897 in Melbourne
 - : large historical isolate collection back to 1950s
- :: National reference laboratory
 - : Salmonella, Listeria, EHEC
- :: WHO regional reference lab
 - : vaccine preventable invasive bacterial pathogens

Bacterial genomics

Small genome







6,000,000,000 letters

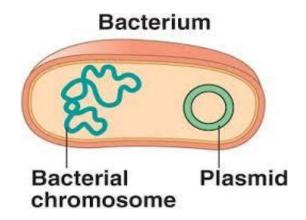
Genome A T G C

3,000,000 letters

30,000 genes

es 3,000 genes

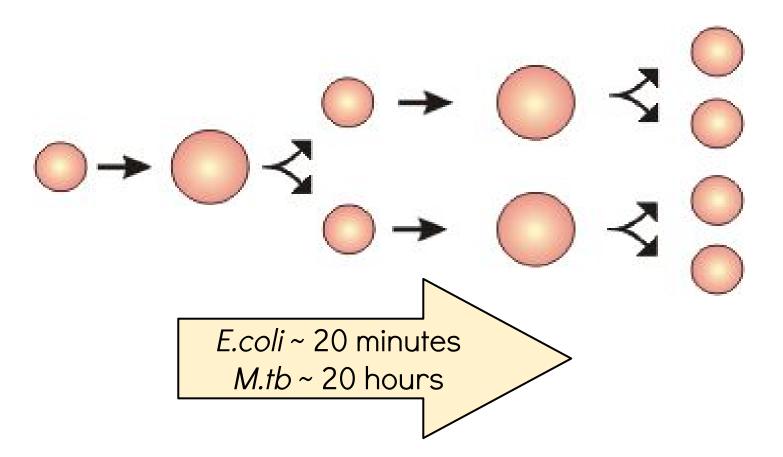
Replicons



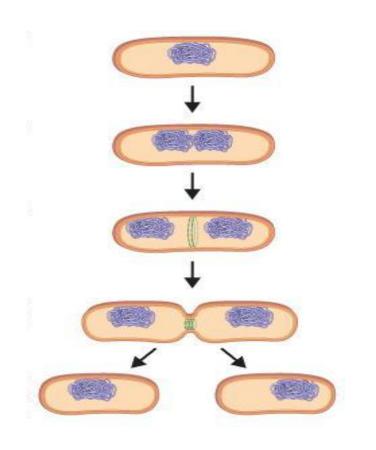
Usually 1 big circular chromosome (1M to 10M bases)

Sometimes 1-6
"mini" chromosomes
(4k - 300k bases)

(Relatively) fast growers



Vertical transfer of DNA



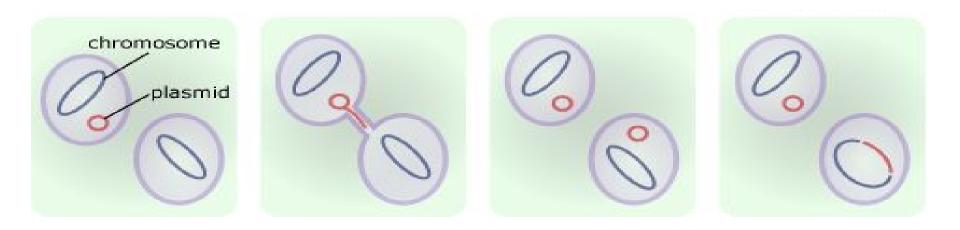
Occurs during cell division

Sometimes it makes an error copying the DNA

 $eg. A \rightarrow T$

Horizontal transfer of DNA

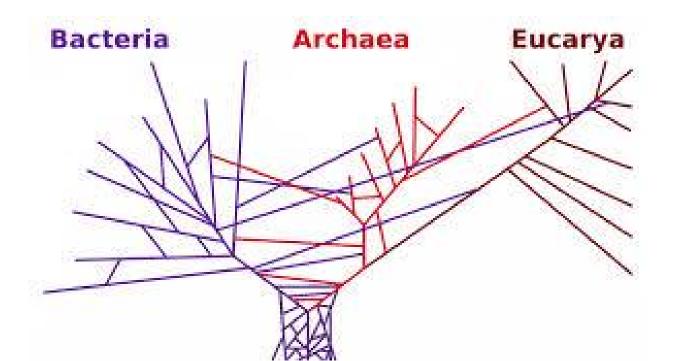
Occurs between bacterial cells



Plasmids: virulence & antibiotic resistance genes!

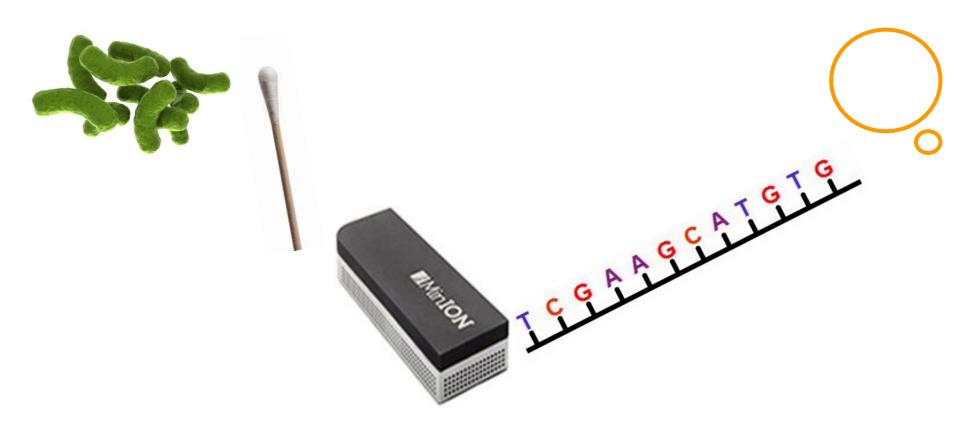
The real tree of life

A mixture of horizontal & vertical transmission.

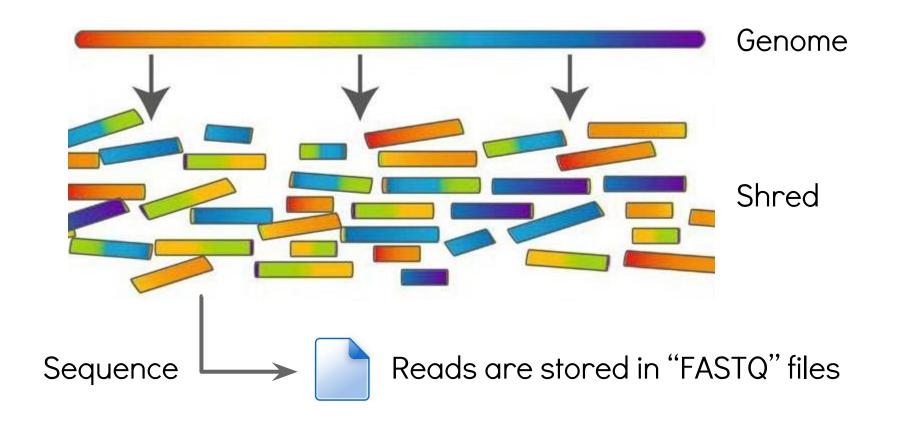


Whole genome sequencing (WGS)

In an ideal world



The real world (for now)



WGS technologies

illumına

100 - 300 bp



100 - 400 bp

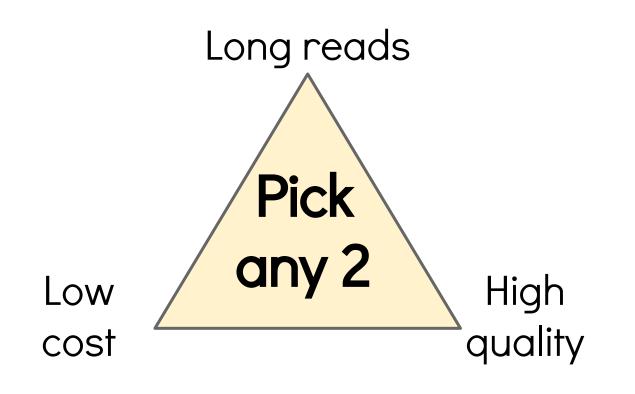


5,000 - 25,000+ bp



5,000 - 150,000+ bp

Which sequencing platform?





Workhorse for bacterial sequencing

illumına

100 - 300 bp



100 - 400 bp

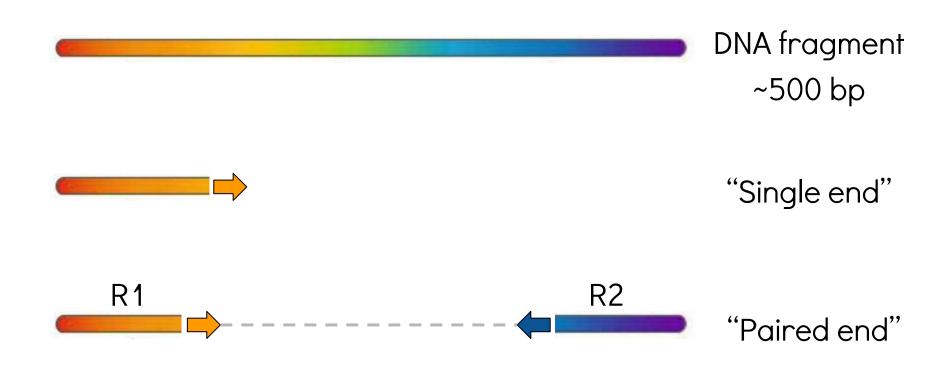


5,000 - 25,000+ bp



5,000 - 150,000+ bp

Types of reads



What you get back

Millions to billions of reads (big files):

... <--- 100 to 300 bp --->

<- 1st read

<- 2nd read

<- last read

Applications of WGS



:: Diagnostics

- : species ⇒ subspecies ⇒ strain identification
- : in silico antibiogram and virulence profile

:: Surveillance

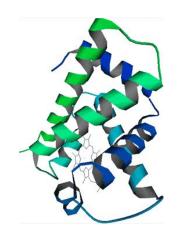
- : in silico genotyping MLST, serotyping, VNTR, MLVA
- : what's lurking in our hospital/community?

- : outbreak detection & source tracking
- : phylogenomics

Not just genomes!

If you can transform your assay into sequencing lots of short pieces of DNA, then NGS is applicable.

- exome (targeted subsets of genomic DNA)
- RNA-Seq (transcripts via cDNA)
- ChIP-Seq (protein:DNA binding sites)
- HITS-CLIP (protein:RNA binding sites)
- methylation (bisulphite treatment of CpG)



FASTA sequence file format

FASTA



FASTA components

Start Sequence description Sequence ID symbol (spaces allowed) (no spaces) >NM 006361.5 Alchohol dehydrogenase (ADH) EC 1.1.1.1 ATGTGCGTCAAGACGGCCGTGCTGAGCGAATGCAGGCGACTTGCGAGCTGGGAGCGAT TTGGATTCCCCCGGCCTGGGTGGGGAGAGCGAGCTGGGTGCCCCCTAGATTCCCCGCC GGATATCTGGGAGCGGGAGGGGGGGGGAATTGA The sequence (usually 60 letters per line)

Multi-FASTA



Concatenation of individual FASTA entries, using ">" as an entry separator

>read00001

TCTTGCGTCAAGACGGCCGTGCTGAGCGAATGCAGGCGACTTGCGAGCTGGGAGCGA

>read00002

TGGATTCCCCCGGCCTGGGTGGGGAGAGCGAGCTGGGTGCCCCCTAGATTCCCCGCC

>read00003

>read00004

TCTGGGAGCGGGGGGGGGGGAATCTGGAGCGAGCTGGGTGCCCCCTAGATTCCCC

>read00004

GCGGAATCTGGAGCGAGCTGGGTGCCCCCTAGATTCCCCGCATCGTAGATTAGATAT

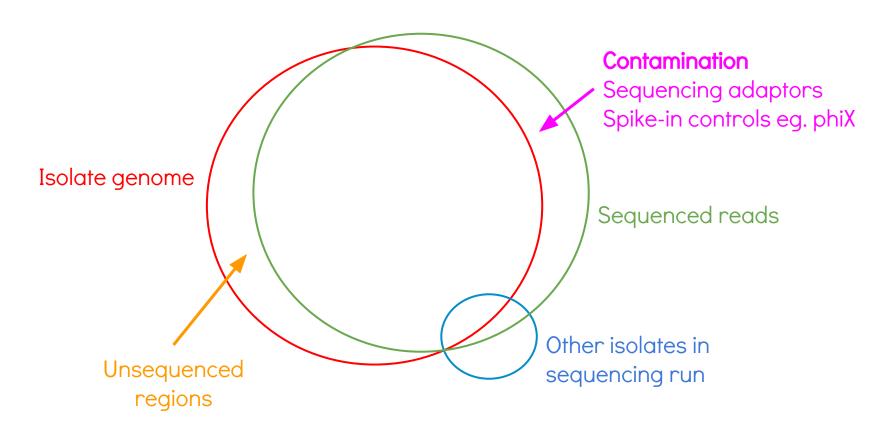
The DNA alphabet

- Standard
 - o AGTC
- Extended
 - adds N (unknown base)
- Full
 - adds R Y M S W K V H D B (ambiguous bases)
 - R = A or G (puRine)
 - \circ Y = C or T (pYrimidine)
 - ... and so on for all the combinations

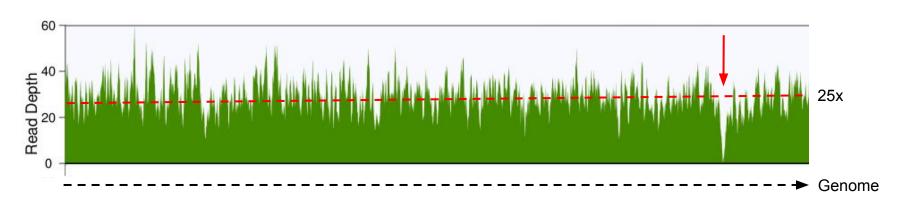


Sequence Quality

What data do we really have?



Do we have enough data?



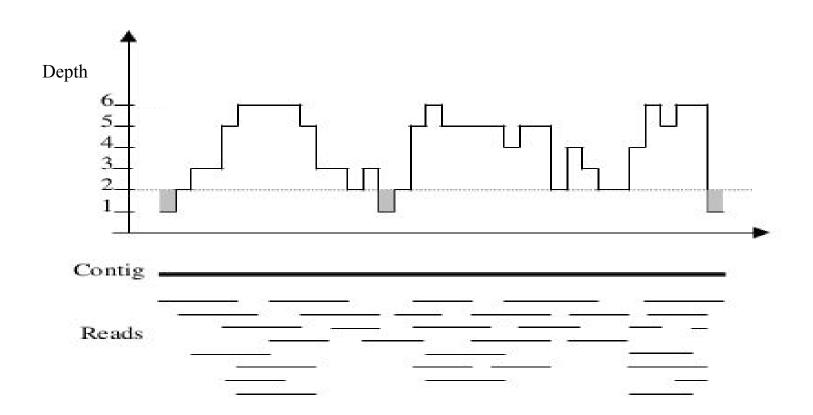
∷ Depth

- : expressed as fold-coverage of genome eg. 25x
- : means each base sequenced 25 times (on average)

:: Coverage

: the % of genome sequenced with depth > 0

Depth (of coverage)



Sequences have errors

nonsense reads

duplicate reads

adaptor read-through

indel errors

uncalled base

substitution errors

instrument oddness

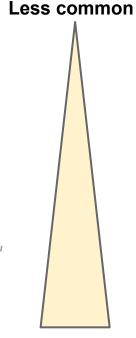
amplify a low complexity library

fragment too short

skipping bases, inserting extra bases

couldn't reliably estimate, replace with "N"

reading wrong base



More common

Illumina reads

Usually 100 - 300 bp

• Indel errors are rare

- Substitution errors < 1%
 - Error rate higher at 3' end
- Very high quality overall



DNA base quality



DNA sequences often have a quality value associated with each nucleotide

- A measure of reliability for each base
- Derived from a physical process
 - chromatogram (Sanger sequencing)
 - pH reading (Ion Torrent sequencing)
 - Voltage measurement (Oxford Nanopore)



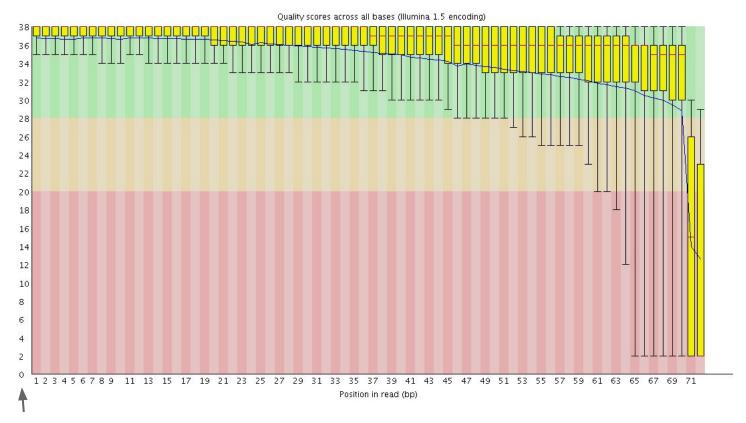


Quality	Chance it's wrong	Accuracy	Description
10	1 in 10	90%	Bad
20	1 in 100	99%	Maybe
30	1 in 1000	99.9%	OK
40	1 in 10,000	99.99%	Very good
50	1 in 100,000	99.999%	Excellent

$$Q = -10 \log_{10} P$$
 <=> $P = 10^{-Q/10}$

Q = Phred quality score P = probability of base being incorrect

Quality plot (FastQC)



Y-axis is "Phred" quality values (higher is better)

Quality filtering

- Keep all reads
 - let the downstream software cope
- Reject some reads
 - average quality below some threshold
 - contain any ambiguous bases
- Trim reads
 - remove low quality bases from end(s)
 - keep longest "sub-read" that is acceptable
- Best strategy is analysis dependent

FASTQ files

FASTQ

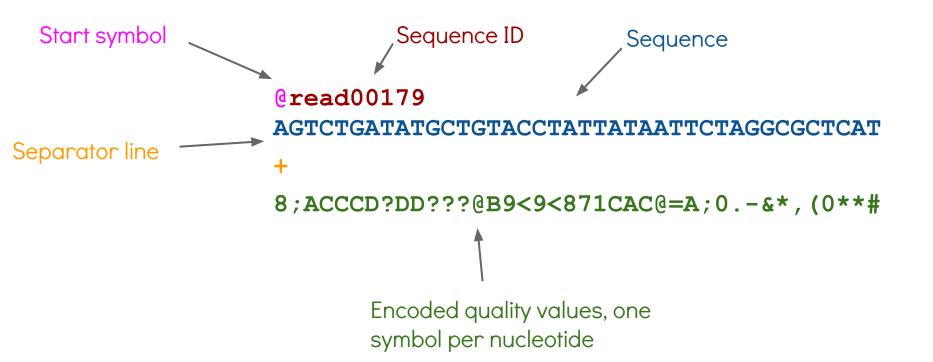


FASTQ sequence entry looks like this:

```
@read00179
AGTCTGATATGCTGTACCTATTATAATTCTAGGCGCTCAT
+
8;ACCCD?DD???@B9<9<871CAC@=A;0.-&*,(0**#</pre>
```

FASTQ components





FASTQ quality encoding



Uses letters/symbols to represent numbers:

Multi-FASTQ

Same as multi-FASTA, just concatenate:

@M00267:3:15997:1501

CTCGTGCTCTACTTTAGAAGCTAATGATTCTGTTTGTAGAACATTTTCTACCACTACATCTTTTTCTTGCTTCGCATCTT

+

@M00267:3:15997:1505

GCCTATAGTAGAAGAAGAAGAAGTGGCTCAAGAAATGAGTGCACCGCAGGAAGTTCCAGCGGCTGAATTACTTCATGAAA

+

<@@FFF?DHFHGHIIIFGIIGIGICDGEGCHIIIIIIIIIIIIIIFG<DA7=BHHGGIEHDBEBA@CECDD@CC>CCCAC

@M00267:3:14073:1508

GTCTTGCTAAATTAAATAATCTGAAATAATTTGTTCTGCCCGGTCCAATTCAGCTAATACGAGACGCATATAATCCTTA

+

@M00267:3:14073:1513

ACGTACAGAGATGCAAAAGTCAGAGAAACTTAATATTGTAAGTGAGTTAGCAGCAAGTGTTGCACATGAGGTTCGAAATC

+

100DDADHGDF?FBGGAFHHCHGGCGGFHIECHGIIGIGFGHGHIIHHEGCCFCB>GEDF=FCFBGGGD0HEHE9=; AD

Got my reads, now what?



Two options



- :: De novo genome assembly
 - : reconstruct original sequence from reads alone
 - : like a giant jigsaw puzzle
 - : Create

- :: Align to reference
 - : identify where each read fits on a related genome
 - : can not always be uniquely placed
 - : Compare

Genome assembly

(the red pill)

De novo genome assembly

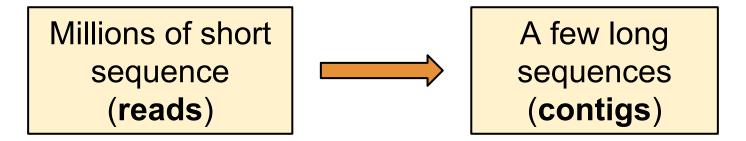
De novo = without reference to other genomes



De novo assembly

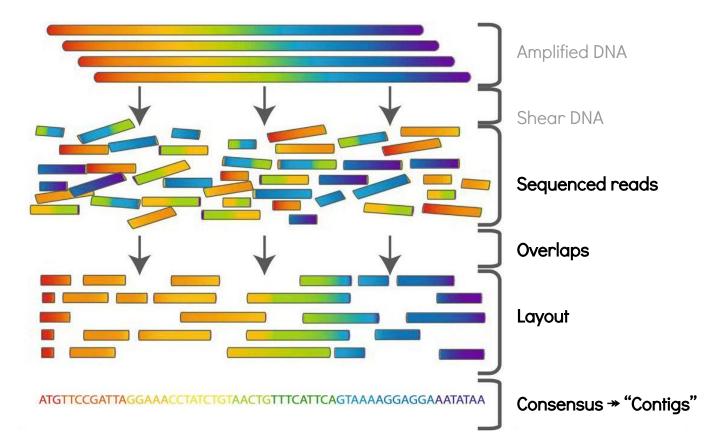


Reconstruct the original genome sequence from the sequence reads only

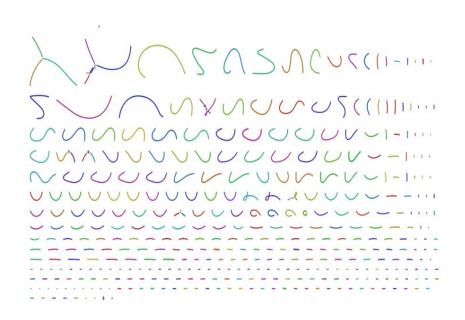


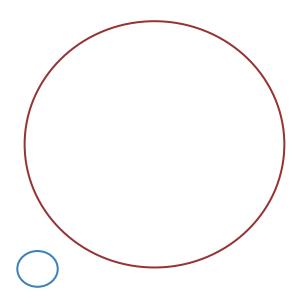
Ideally, one sequence per chromosome.

Overlap - Layout - Consensus



Draft vs Finished genomes

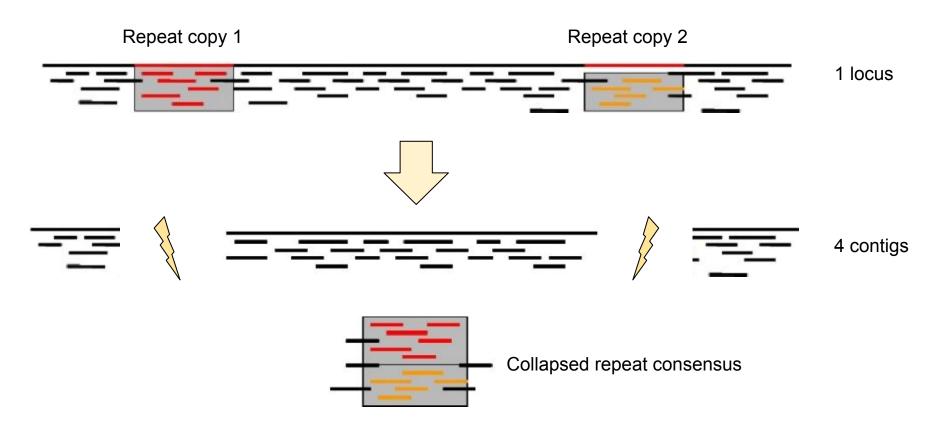




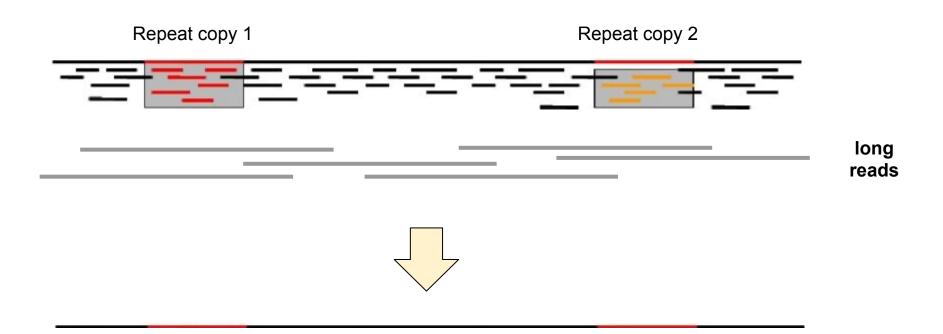
250 bp - Illumina - \$250

8000 bp - Pacbio - \$2500

Repeats



Long reads can span repeats



The law of repeats



 It is impossible to resolve repeats of length L unless you have reads longer than L.

 It is impossible to resolve repeats of length L unless you have reads longer than L.

Align to reference

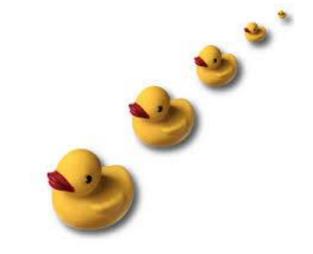
(the blue pill)

Align to reference

Seven short 4bp reads

AGTC TTAC GGGA CTTT

TAGG TTTA ATAG



Aligned to 31bp reference

AGTCTTTATTATAGGGAGCCATAGCTTTACA

AGTC

TAGG

ATAG

TTAC

TTTA

GGGA

CTTT

Ambiguous alignment

Eight short 4bp reads

AGTC TTAC GGGA CTTT

TAGG TTTA ATAG TTAT



Aligned to 31bp reference

AGTCTTTATTATAGGGAGCCATAGCTTTACA

AGTC TAGG ATAG TTAC

TTTAT GGGA CTTT

TTAT

D'oh!

Look at differences

AGTCTGATTAGCTTAGCTTGTAGCGCTATATTAT

Reference

AGTCTGATTAGCTTAGAT

ATTAGCTTAGATTGTAG

CTTAGATTGTAGC-C

TGATTAGCTTAGATTGTAGC-CTATAT

TAGCTTAGATTGTAGC-CTATATT

TAGATTGTAGC-CTATATTA

TAGATTGTAGC-CTATATTAT

Reads

Look at differences

Substitution

Deletion

AGTCTGATTAGCTTAGCTTGTAGCGCTATATTAT

Reference

AGTCTGATTAGCTTAGAT

ATTAGCTTAGATTGTAG

CTTAGATTGTAGC-C

TGATTAGCTTAGATTGTAGC-CTATAT

Reads

TAGCTTAGATTGTAGC-CTATATT

TAGATTGTAGC-CTATATTA

TAGATTGTAGC-CTATATTAT

Types of variants we can detect

Туре	Reference	Alternate
SNP (single)	Т	G
MNP (multiple)	TA	GC
Insertion	AGT	ACGT
"	ATCGGG	ATC TGA GGG
Deletion	ACGT	AGT
"	ATC <u>TGA</u> GGG	ATCGGG

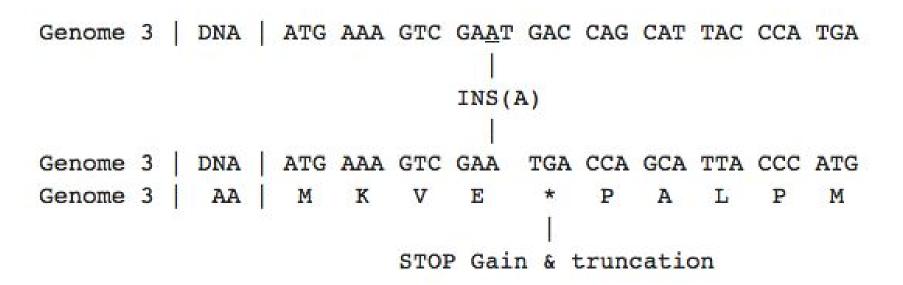
Synonymous & non-synonymous SNPs

```
ATG AAA GTT GAT GAC CAG CAT TCC CCA TGA
                 ATG AAA GTC GAT GAC CAG CAT TAC CCA TGA
                        SNP(T=>C)
                                            SNP(C=>A)
Genome 1
            AA
Genome 2
            AA
                          SYN
                                             NON-SYN
```

Indel causing frame-shift

```
ATG AAA GTT GAT GAC CAG CAT TCC CCA TGA
Genome 1
                  M
                      K
                                          H
                 ATG AAA GTC -AT GAC CAG CAT TAC CCA TGA
                           DEL(G)
                 ATG AAA GTC ATG ACC AGC
                                         ATT ACC CAT GA?
Genome 2
            AA
                                             STOP Loss & read-through
```

Indel causing truncation



Best practice

- Use <u>both</u> approaches
 - □ reference-based + *de novo*



- Best of both worlds
 - $\ \square$ and worst of both worlds interpretation is <u>non-trivial</u>
- Still need
 - good epidemiology, metadata and domain knowledge!

Downstream bioinformatics

Correcting contigs



- Pacbio long reads
 - □ ~10 x 1bp homopolymer insertion errors per Mbp
 - □ Cause frame-shift errors in coding genes
- Illumina short reads
 - Don't suffer from indel issues
- Use Illumina to correct Pacbio!
 - □ It works well

Species identification

Name	Readco	Readcount (% of classified reads)	
	Mycobacterium abscessus	959 (11.65%)	
	Vibrio cholerae	937 (11.38%)	
	Salmonella enterica subsp. enterica	738 (8.96%)	
	Enterobacter cloacae subsp. cloacae NCTC 9394	704 (8.55%)	
	Rhodobacter sphaeroides 2.4.1	613 (7.44%)	
	Staphylococcus aureus	567 (6.89%)	
	Klebsiella	563 (6.84%)	
	Bacillus cereus ATCC 10987	434 (5.27%)	
	Bacillus cereus	347 (4.21%)	
	Rhodobacter sphaeroides	334 (4.06%)	
	(Remaining organisms)	2038 (24.75%)	



Annotation

Adding biological information to sequences.

ribosome binding site

delta toxin
PubMed: 15353161

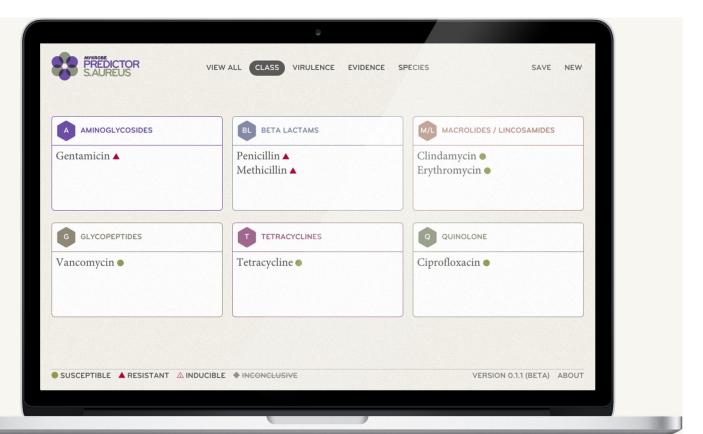
tandem repeat

transfer RNA

Leu-(UUR)

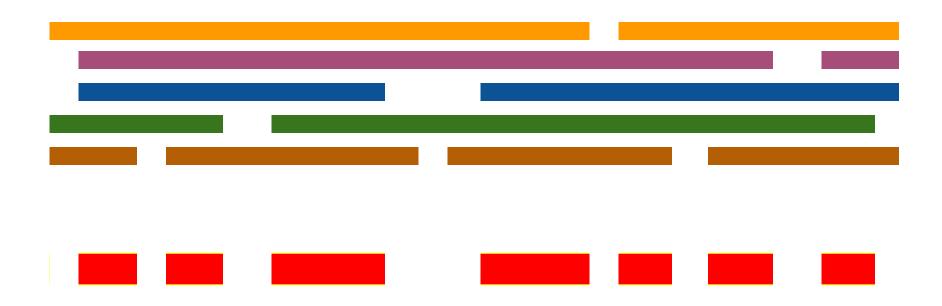
homopolymer 10 x T

Antimicrobial resistance



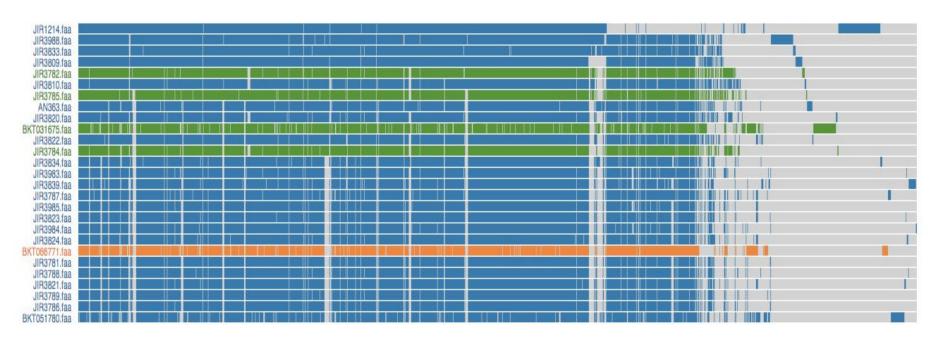
Phylogenetics

Core genome



Core is <u>common to all</u> & has <u>similar</u> sequence.

Pan genome



Rows are genomes, columns are genes.

Application to public health and clinical microbiology

Traditional workflow



A bacterial isolate



Focus on a small "informative" section



e.g. MLST, VNTR, PFGE, <insert genotyping method here>

Genotype shows isolates are related

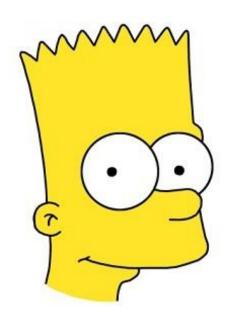






D'oh!

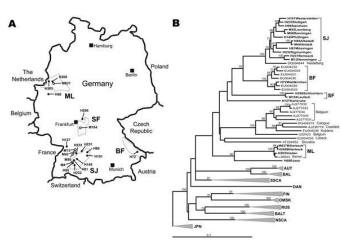


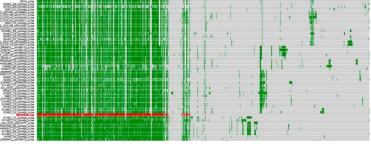




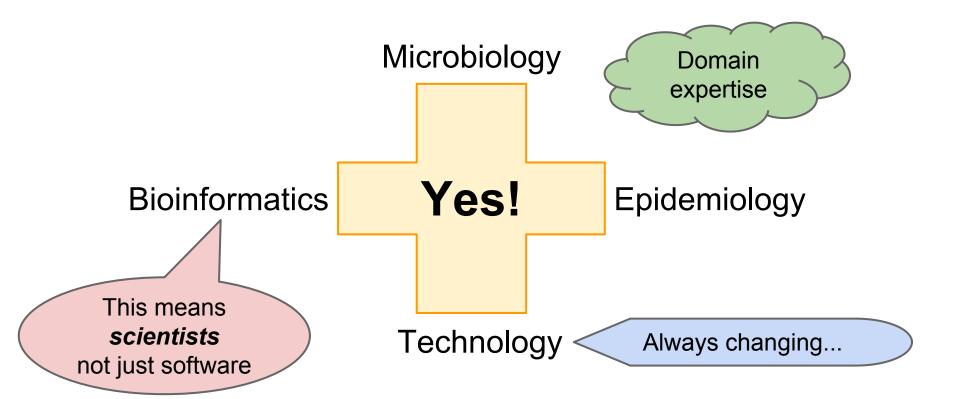
Modern workflow







Does WGS deliver?



Conclusions

Take home messages

- Shorts reads & repeats
 - □ Can't reconstruct genome fully
 - □ Can't always align unambiguously
- Bioinformatics is a skill
 - not just pushing buttons
 - be nice to us and we will enjoy helping you

Acknowledgements

Anne Syme

Marcel Behr

Simon Gladman

Lynn Dery Capes

Anders da Silva

Robyn Lee

Google

Ines Levade

Fin. Merci.

That's all folks!