

Comparative Genomics and Visualisation

BS32010

2. Bulk Genome Properties



**The James
Hutton
Institute**

Leighton Pritchard^{1,2,3}

¹Information and Computational Sciences,

²Centre for Human and Animal Pathogens in the Environment,

³Dundee Effector Consortium,

The James Hutton Institute, Invergowrie, Dundee, Scotland, DD2 5DA



Acceptable Use Policy

Recording of this talk, taking photos, discussing the content using email, Twitter, blogs, etc. is permitted (and encouraged), providing distraction to others is minimised.

These slides will be made available on SlideShare.

These slides, and supporting material including exercises, are available at https://github.com/widdowquinn/Teaching-2015-03-17-UoD_compgenvis



Table of Contents

Bulk Genome Comparisons

Experimental bulk genome comparisons

The impact of high throughput sequencing

In silico bulk genome comparisons



Bulk genome comparisons

You don't have to sequence genomes to
compare them
(but it helps)
PART 1: Experimental



Genome comparisons predate NGS

- Sequence data wasn't always cheap and abundant
- Practical, experimental genome comparisons were needed





Bulk genome comparisons

**Calculate values for individual genomes,
then compare them.**



Bulk genome properties

- Large-scale summary measurements
- Measure genomes independently - compare values later



Bulk genome properties

- Large-scale summary measurements
- Measure genomes independently - compare values later
 - What kinds of measurements/properties?



Bulk genome properties

- Large-scale summary measurements
- Measure genomes independently - compare values later
 - Number of chromosomes
 - Ploidy
 - Chromosome size
 - Nucleotide (A,C,G,T) frequency



Chromosome stains ^a ^b

^aIRGSP (2005) *Nature* doi:10.1038/nature03895

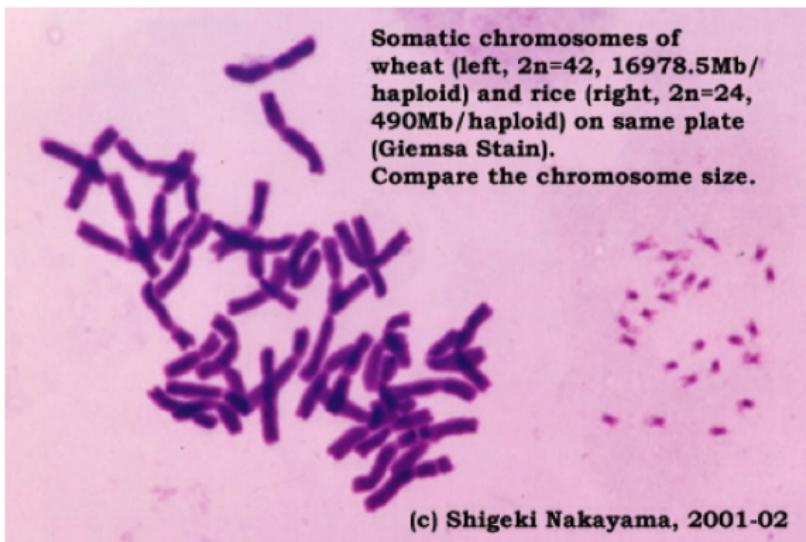
^bBrenchley *et al.* (2012) *Nature* doi:10.1038/nature11650



The James
Hutton
Institute

Count chromosomes, estimate size

Wheat: 17Gbp; Rice: 390Mbp





Chromosome count/size ^a ^b

^a Kamisugi et al. (1993) *Chromosome Res.* 1(3): 189-196

^b Wang et al. (2013) *Nat. Rev. Genet.* doi:10.1038/nrg3375

- chromosome count and ploidy can vary widely

Organism	Chromosomes	Ploidy
<i>E. coli</i>	1	1
Human (<i>H. sapiens</i>)	46	2
Rice (<i>O. sativa</i>)	24	1
Adders-tongue (<i>Ophioglossum reticulatum</i>)	1260	84
Domestic (not wild) wheat somatic	42	6
Domestic (not wild) wheat gametic	14	2

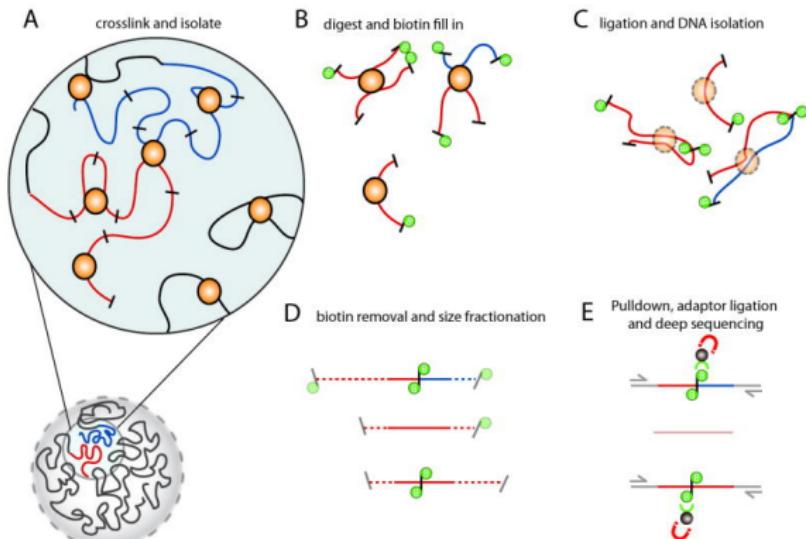




Physical study of chromosomes ^a

^a Belton et al. (2012) *Methods* doi:10.1016/j.ymeth.2012.05.001

- Genome size and chromosome count do **not** indicate organism “complexity”
- Modern physical study of chromosomes still produces surprises (e.g. Hi-C: chromatin interaction)





Nucleotide content ^a

^a Krane et al. (1991) *Nucl. Acids Res.* 19(19): 5181-5185

- Radiolabel monophosphates from genomic DNA
- Separate by thin-layer chromatography
- Compare label ratios using scanner

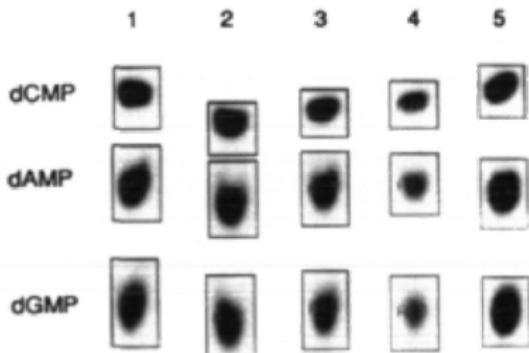


Table 4. Base composition of yeast artificial chromosomes containing of human DNA inserts.

	YAC	Replicate	Corrected		%G-C	%G-C	Average	GC content	Std. dev.
			GA	CA	from GA	from CA			
dCMP	yKM19-3 (180 kb)	1	0.737	0.799	42.4	44.4	43.4		
		2	0.703	0.771	41.3	43.5	42.4		0.58
		3	0.707	0.768	41.4	43.4	42.4		
dAMP	yCF-4 (330 kb)	1	0.724	0.768	42.0	43.4	42.7		
		2	0.708	0.731	41.4	42.2	41.8		0.45
		3	0.701	0.759	41.2	43.2	42.2		
dGMP	yW30-5 (300 kb)	1	0.722	0.747	41.9	42.7	42.3		
		2	0.744	0.765	42.7	43.3	43.0		0.47
		3	0.713	0.739	41.6	42.5	42.1		
	yJ311-2 (230 kb)	1	0.725	0.747	42.0	42.8	42.4		
		2	0.728	0.757	42.1	42.4	42.6		0.10
		3	0.728	0.738	42.1	42.4	42.5		
	yJ311-5 (200 kb)	1	0.695	0.726	41.0	42.1	41.5		
		2	0.746	0.756	42.7	43.0	42.9		0.70
		3	0.731	0.733	42.2	42.3	42.3		
	yHPRT (680 kb)	1	0.718	0.744	41.8	42.7	42.2		
		2	0.718	0.725	41.8	42.0	41.9		0.46
		3	0.694	0.713	41.0	41.6	41.3		



Table of Contents

Bulk Genome Comparisons

Experimental bulk genome comparisons

The impact of high throughput sequencing

In silico bulk genome comparisons



This happened...

- Cheap, accurate, high-throughput sequencing





Four different chemistries ^a

^aLoman *et al.* (2012) *Nat. Rev. Micro.* **31**:294-296 doi:10.1038/nbt.2522



Reads differ by technology, and can require different bioinformatic treatment...

- **Roche/454:** Pyrosequencing (long reads, but expensive, and high homopolymer errors) (700-800bp, 0.7Gbp, 23h)
- **Illumina:** Reversible terminator (cost-effective, massive throughput, but short read lengths) (2x150bp, 1.5Gbp, 27h)
- **Ion Torrent:** Proton detection (short run times, good throughput, high homopolymers errors) (200bp, 1Gbp, 3h)
- **PacBio:** Real-time sequencing (very long reads, high error rate, expensive) (3-15kbp, 3Gbp/day, 20min)

... different error profiles, varying capability to assemble/determine variation



Costs of sequencing^a

^a Miyamoto et al. (2014) *BMC Genomics* 15:699 doi:10.1186/1471-2164-15-699

Cost and required DNA comparison

	GS Jr	Ion PGM	MiSeq	PacBio
Instrument cost	\$108K	\$50K	\$99K	\$900K
Sequence yield per run	35Mb	2Gb (400bp)	8 Gb	1 Gb/8 SMRT cells
Running cost	\$1000/1run(35Mb)	\$437/Gb	\$93/Gb	\$1800/Gb
Sequence Run time	10 hr	7.3 hr	39 hr	16 hr/8 SMRT cells
Other time consuming steps	Library prep: 3hr emPCR: 6hr	Library prep: 3.5 hr emPCR : 8 hr	Library prep: 7 hr	Library prep: 5 hr
DNA requirements	500 ng with 1.8 OD	250 ng	250 ng	100ng (250bp library) - 5μg (20kb library)



After that, the flood...

High-throughput sequencing methods have completely changed the landscape of biology

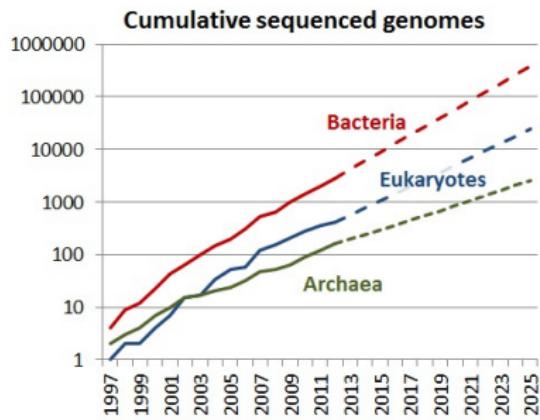
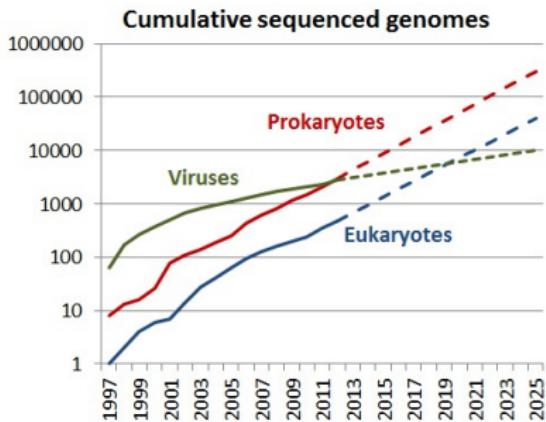
(Nearly) complete, (mainly) accurate sequence data is now inexpensive (and cheaper than analysis)

- GOLD (19/2/2014): 3,011 “finished” ; 9,891 “permanent draft” genomes
- GOLD (18/11/2014): 6,649 “finished” ; 23,552 “permanent draft” genomes
- NCBI WGS (19/2/2014): 17,023 microbial genomes
- NCBI WGS (18/11/2014): 26,026 microbial genomes



Predicting the future is hard...

Su *et al.* attempted to answer this¹:

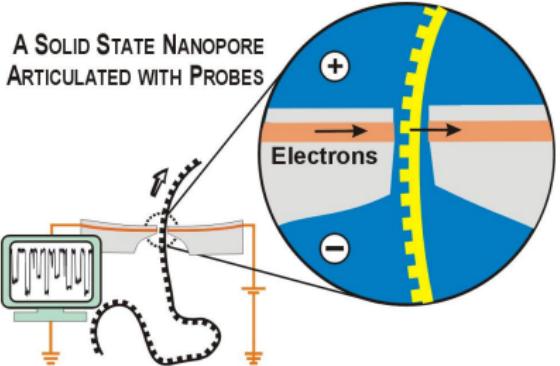
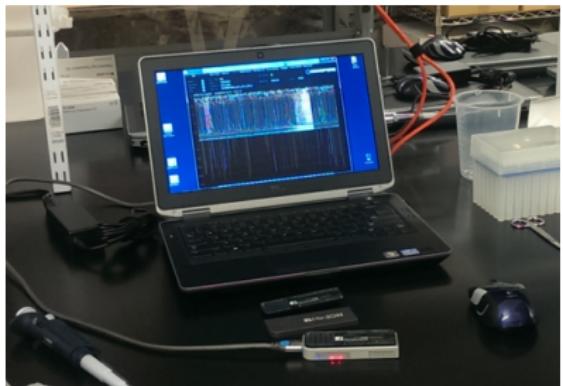


¹ <http://sulab.org/2013/06/sequenced-genomes-per-year/>



What's coming next?

Oxford Nanopore. A sequencer the size of your hand.



- Microfluidics, single-molecule sequencing; 11-70kbp reads
- Reports current across pore (tiny electron microscope) as molecule moves through
- \$10/Mbp, 110Mbp per flowcell²

²Yaniv Erlich (2013) Future Continuous blog



Table of Contents

Bulk Genome Comparisons

Experimental bulk genome comparisons

The impact of high throughput sequencing

In silico bulk genome comparisons



Bulk genome comparisons

You don't have to sequence genomes to
compare them
(but it helps)
PART 2: in silico



Bulk genome comparisons

EXERCISE 1:

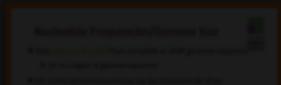
`ex01_gc_content.ipynb`



Nucleotide frequency/genome size

Very easy to calculate from complete or draft genome sequence

```
In [1]: from Bio import SeqIO
In [2]: s = SeqIO.read("data/NC_000912.fna", "fasta")
In [3]: a, c, g, t = s.seq.count("A"), s.seq.count("C"), s.seq.count("G"), s.seq.count("T")
In [4]: float(g + c)/len(s)
Out[4]: 0.40008010837904245
In [5]: float(g - c)/(g+c)
Out[5]: 0.002397259225467894
```



GC content, chromosome size can be characteristic of an organism.



Blobology ^{a b}

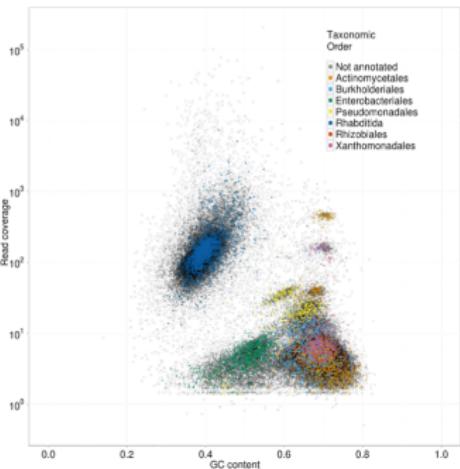
^aKumar & Blaxter (2011) *Symbiosis* doi:10.1007/s13199-012-0154-6

^b<http://nematodes.org/bioinformatics/blobology/>

Sequence data can be contaminated by other organisms

- Host and symbiont DNA have different %GC
- Host and symbiont DNA differ in coverage

- Assemble genome
- Map reads
- Plot coverage against %GC





Nucleotide k -mers

Sequence data is necessary to determine k -mers/frequencies
Not possible by experiment

- Nucleotides, $k = 1$, 4 \times 1-mers
A, C, G, T
- Dinucleotides, $k = 2$, 16 \times 2-mers
AA, AC, AG, AT, CA, CC, CG, CT, GA, GC, GG, GT, TA, TC, TG, TT
- Triucleotides, $k = 3$, 64 \times 3-mers
- k -nucleotides, $4^k \times k$ -mers



Finding k -mers

ACTIVITY:

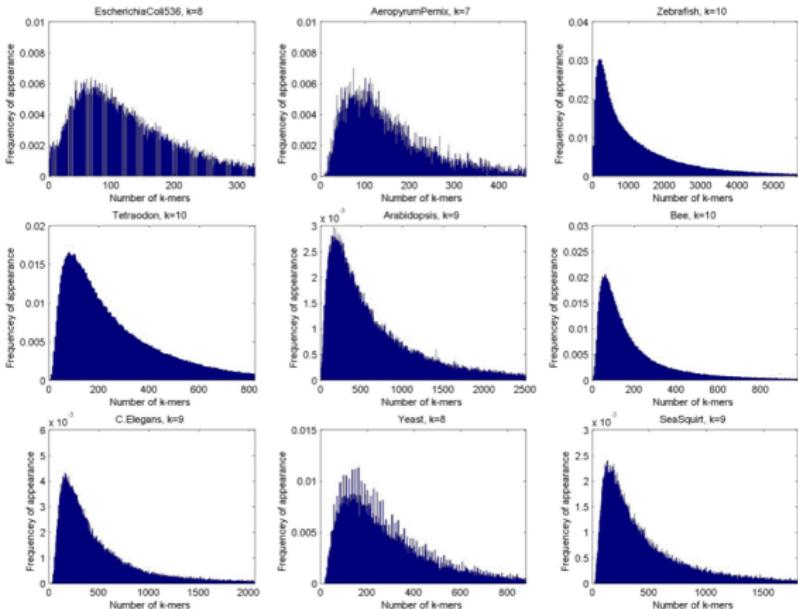
https://widdowquinn.shinyapps.io/nucleotide_frequencies/



k-mer spectra ^a

^aChor et al. (2009) *Genome Biol.* doi:10.1186/gb-2009-10-10-r108

k-mer spectrum: frequency distribution of observed *k*-mer counts.
Most species have a unimodal *k*-mer spectrum



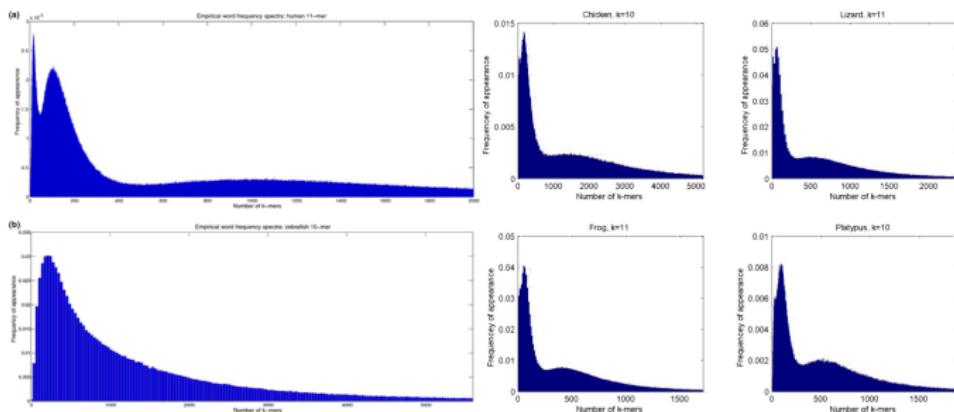


k-mer spectra ^a

^aChor et al. (2009) *Genome Biol.* doi:10.1186/gb-2009-10-10-r108

All mammals tested (and some other) species have *multimodal k-mer spectra*

Genomic regions differ in this property





Bulk genome comparisons

EXERCISE 2:

`ex02_kmer_spectra.ipynb`



Licence: CC-BY-SA

By: Leighton Pritchard

This presentation is licensed under the Creative Commons Attribution ShareAlike license

<https://creativecommons.org/licenses/by-sa/4.0/>