

Introduction

Isheng Jason Tsai

Introduction to NGS Data and Analysis
Lecture 1



Welcome

Lecture objective

Introduction

Some basics

Sequencing platforms

Data types

Analysis

This course is called “Introduction to Next-Generation Sequencing (NGS) Data and Analysis”

Actually

- Next Generation Sequencing is really “now” sequencing
- It won’t be so easy to tell you everything about NGS (it’s a bit like saying what can we do with PCR?)

About myself

- 2015. January. Assistant research fellow (助研究員) at Academia Sinica
- Was at Sanger Institute, UK since the start of “NGS” for four years.
- 26 publications.
(2 Nature, 1 Science, 2 Nature Genetics, 2 PNAS)
- Background: Genetics, bioinformatics, population genetics, assembly
- So I know a bit about genomics and NGS.

What I expect from you

- I will share *all* of my experiences, esp. from someone who come from a background in biology
- Expose you to different *environment* and *thinking*
- Be pragmatic: don't just do sequencing because you have the budget or your boss gave you the data

Please

- *share , be open and network*

For this course

- 60% theories
- 40% Practical
 - You need a laptop with linux environment installed ;
 - Mac is preferred
- Report at week 11 and 18

It's a big world out there

- Read, read, read
- Setup twitter and follow what others are doing

Tweets Tweets & replies Photos & videos

You Retweeted

OfficialSMBE @OfficialSMBE · 23h
MBE latest: Large-effect beneficial synonymous mutations mediate rapid and parallel adaptation in a bacterium dlvr.it/KbShdx

6 4 ...

You Retweeted

OfficialSMBE @OfficialSMBE · 22h
GBE latest: Genome Resequencing Identifies Unique Adaptations of Tibetan Chickens to Hypoxia and High-dose... dlvr.it/KbSvMX

1 1 ...

You Retweeted

Justin Fay @justinfay · 19h
Check out our paper on *S. paradoxus* in Slovenian vineyards, including our first #vineyard #microbiome journal.frontiersin.org/article/10.338...

8 9 ...

You Retweeted

Rob Waterhouse @rmwaterhouse · Feb 20
Trait databases, data quality, trees, genome structures, disease, biodiversity, @erichjarvis Ann.Rev. #birdgenomes

Erich Jarvis @erichjarvis
My perspective on questions that can be answered when all vertebrate genomes are sequenced @Genome10K @B10K_Project jarvislab.net/wp-content/upl...

1 1 ...

You Retweeted

Sujai @sujaik · Feb 20
For anyone following the ridiculousness in India, this is brilliant scroll.in/article/803856... @Sanjana2808 @karunanundy

1 1 ... View summary

You Retweeted

James Wasmuth @jdwasmuth · Feb 19
Using #PacBio to gain a high-resolution phylogenetic microbial community profile bit.ly/1oR4qde

3 1 ...

Resources

Some very useful websites:

- <http://angus.readthedocs.org/en/2015/#>
- http://www.pasteur.fr/~tekaia/BCGA2014/BCGA2014_Prog.html
- <http://evomics.org/>
- <http://molb7621.github.io/workshop/>
- <https://sequencing.qcfail.com>
- <http://schatzlab.cshl.edu/teaching/>

Quick survey

- Total number of students: XX
- Anyone already have a dataset?
- Anyone about to design their own experiment, produce sequences and analyse themselves?
- Assembly?
- RNAseq?
- Anything else?
- Familiarity with Linux environment?
- Programming experiences?

Always start with a question

Problem

Most people doing genomics not actually doing genomics

Posted on July 27, 2015 by jovialscientist

CAMBRIDGE. Most people who claim to be genomics researchers are not actually doing genomics at all, and instead are just sequencing things and calling it genomics, it has been found.

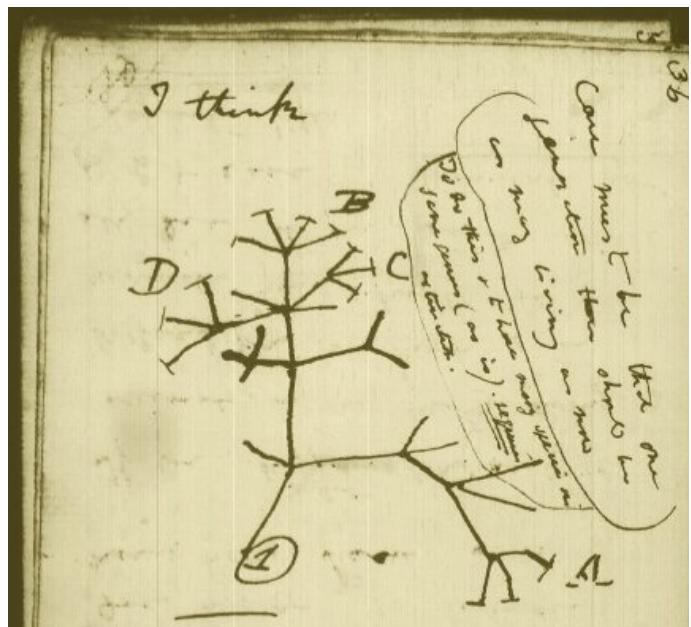
“Genomics is the study of genomes” said Barney Ewingsworth III from the Excellent Biology Institute (EBI) “and genomes are incredibly complex, with repeat regions, duplications, deletions, selective sweeps, gene deserts, 3D structure, mobile elements etc etc. ... and it turns out that many people who say they are genomic researchers are actually just people with a few quid who paid to sequence a stupid genome, like the lesser spotted tree trout. Then they assemble it (badly), submit it to GenBank still full of adapters, and bloody PhiX, and get a paper in *BMC I couldn’t get this into Genome Research*. It’s a scandal – they give genomics a bad name!” he finished, and then went back to his day job as Mayor of London.

In an earlier survey, it was found that many scientists are sequencing things because they can’t think of anything else to do. Now it would appear that those very same scientists have no idea how to handle the data, and are poisoning the well with hundreds of crappy genomes.

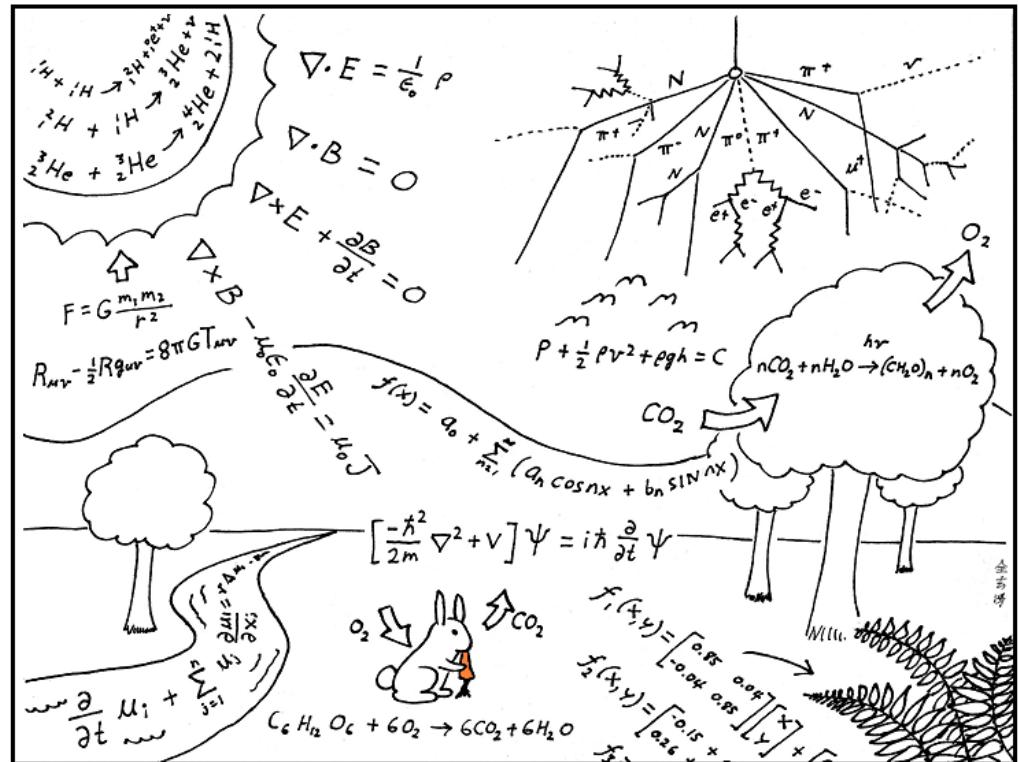
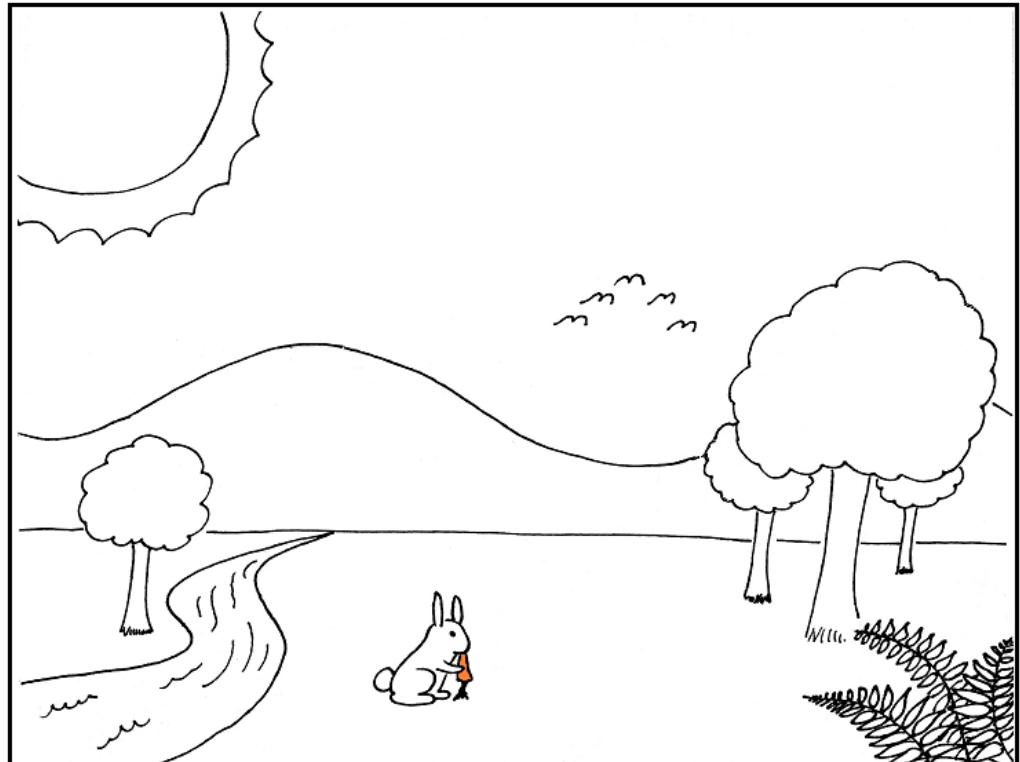
<https://thescienceweb.wordpress.com/2015/07/27/most-people-doing-genomics-not-actually-doing-genomics/>

Nothing makes sense in the light of evolution

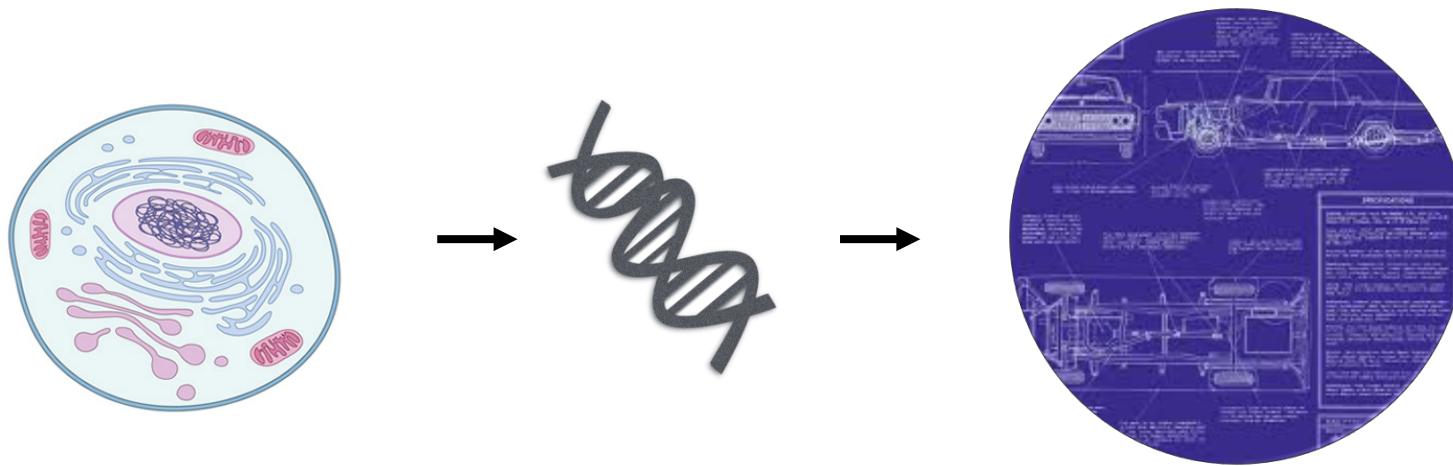
Theodosius Dobzhansky 1973



This is how scientists see the world

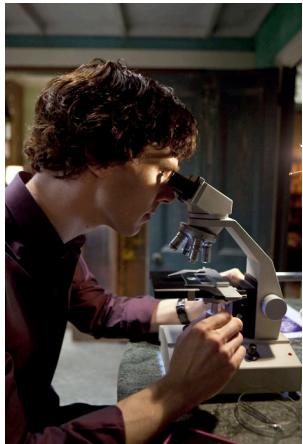


Genome



Genome = Parts list of a single genome

Many perceptions of genomics



What my parents
think I do



What some colleagues
think we do



What more friendly colleagues
think we do



What my friends
think I do



What we
think we do

A black terminal window displaying a dense block of white command-line text, which appears to be a mix of programming code and system logs.

What we
actually do

Why sequence a genome?

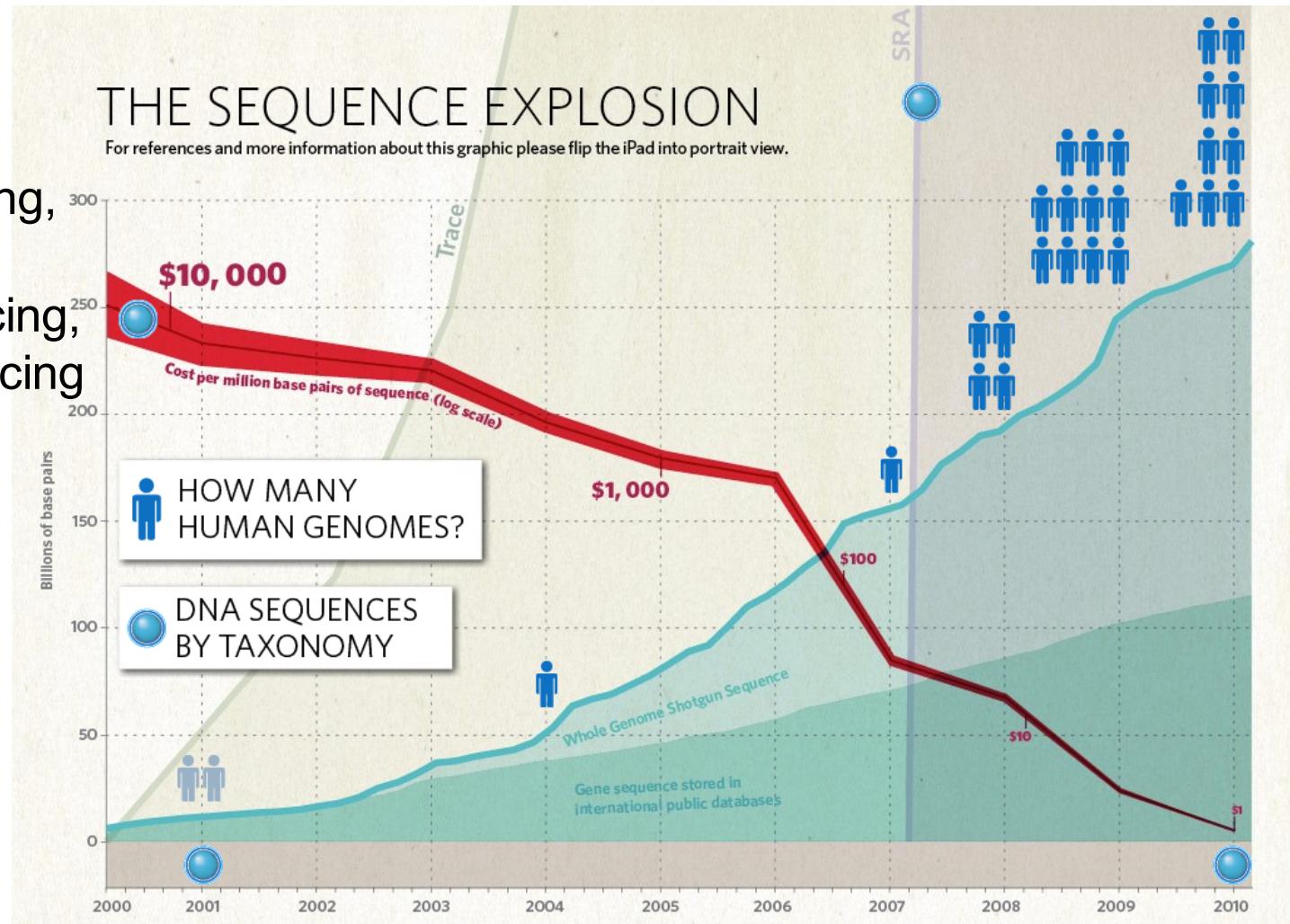
- Phylogenetic position
- Differences between species (comparative genomics)
- Variations between individuals (population genetics)
- Help to understand biology
- Of economic, agricultural, medical, ecology values
- ~~Some lab just had the money ; don't do it~~
- **More case studies: Lecture 2**

NGS dos and don't

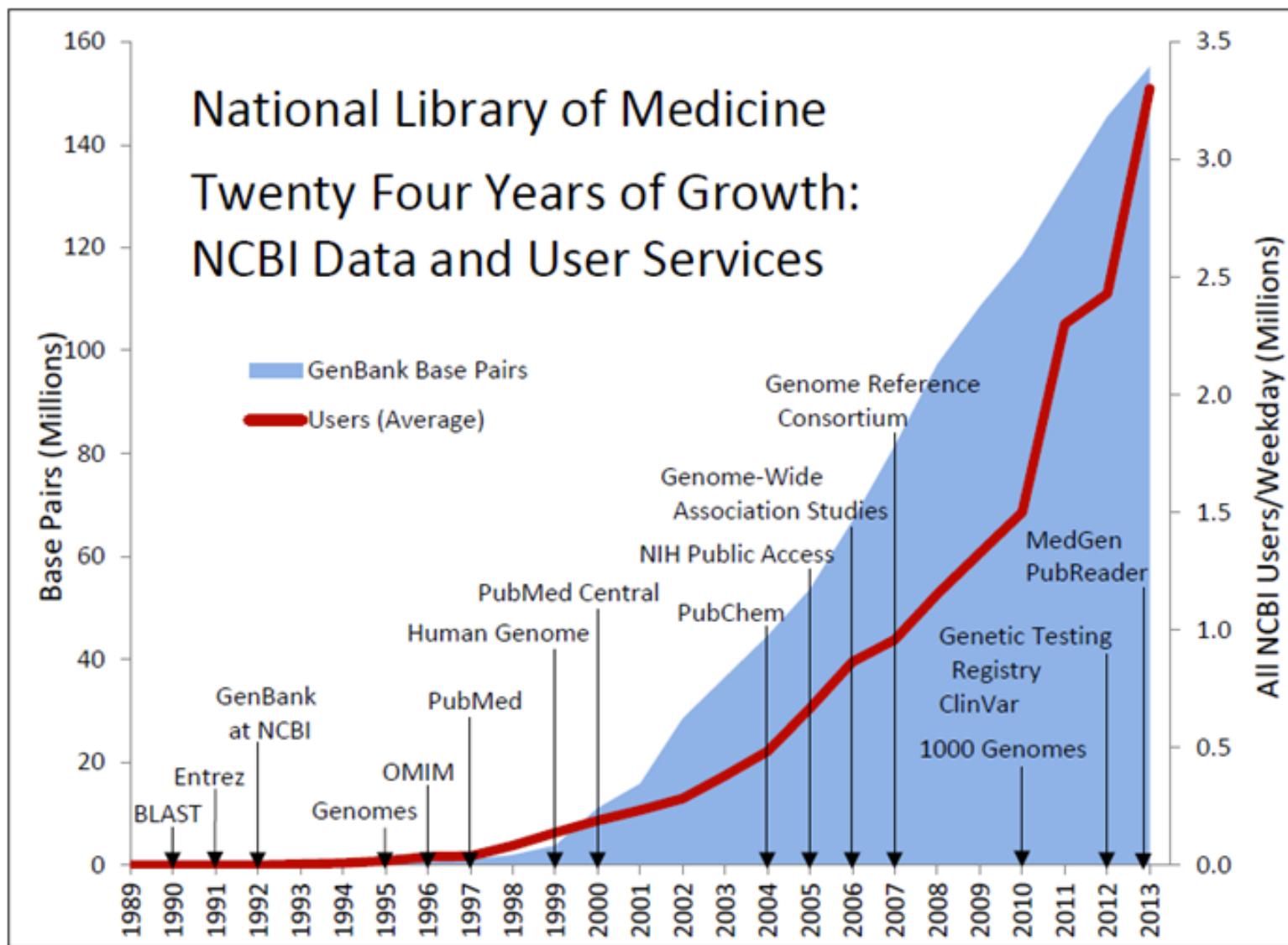
- Embrace it
- Don't just do it without a question
- Don't do it because you can (lots of \$\$, want to jump in)
- Don't just hate it because you don't know how to do it
 - Typical scenario: “We should focus on more traditional methods because NGS is expensive”
 - Typical scenario 2: “These people who do mathematics (?) don't know what ecology/biology/conservation are”
- **More case studies: Lecture 2**

What is NGS?

- = Next generation sequencing,
- = deep sequencing
- = High Throughput Sequencing,
- = Massively parallel sequencing
- = 次世代定序
- = 高速高量定序



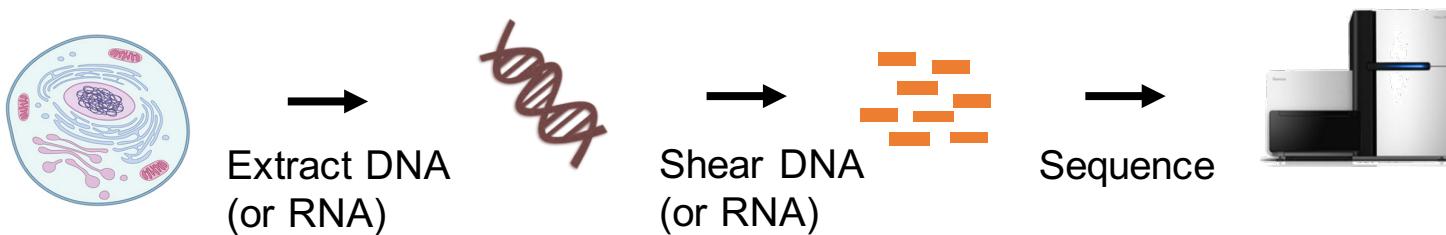
<http://www.nature.com/news/2010/100331/full/464670a.html>



Some basics

A genome project

Wet lab work



Bioinformatics

Data QC
(Lecture 5)



Variant
(Lecture 9,15)

ATCG
AT~~G~~
ATCG

Assembly
(Lecture 4)



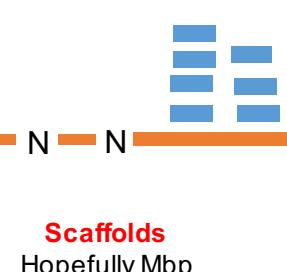
Annotation
(Lecture 8)



Mapping (Lecture 3)
RNAseq (Lecture 6)



Scaffolding
(Lecture 4)



A small project's typical output

Sample Name	Sample ID	Lane ID	Yield (Mb)	# of Reads
F2-1	SG-IB01	1	11,435	75,729,838
F2-2	SG-IB02		12,014	79,561,504
F2-3	SG-IB03		11,577	76,666,714
F3-2	SG-IB05		11,119	73,638,446
F3-4	SG-IB07		10,399	68,870,380
F3-5	SG-IB08		11,671	77,292,976
F3-1	SG-IB09		12,474	82,610,516
F3-3	SG-IB10		11,916	78,915,536
F2-1	SG-IB01	2	11,366	75,271,724
F2-2	SG-IB02		11,920	78,940,010
F2-3	SG-IB03		11,481	76,031,166
F3-2	SG-IB05		11,054	73,203,066
F3-4	SG-IB07		10333	68,429,564
F3-5	SG-IB08		11550	76,488,178
F3-1	SG-IB09		12328	81,640,878
F3-3	SG-IB10		11812	78,225,876

8 exome samples ;

2 Illumina Hiseq lanes with 184GB of data

~100X of human exome to detect disease causing SNP

Higher yield at lower cost = More samples can be barcoded into one lane

More samples = more replicates (power) in statistical analysis to pick up real biological difference

More data but less people with informatics skills

- Sequencing is the result of many types of experiment
- Everyone wants to make use of this technology
- Not everyone will be able analyse them
 - You can't just open the file in Microsoft office anymore
- Collaborate or learn yourself
- **Bottleneck is bioinformatics analysis**

You will end up with an analysis pipeline

Run **multiple programs** to analyse / get the results

Problems:

- Which program to use?
- Which parameter to use for each program?
- How do you get results of program A to feed into program B?
- How do you know if the program finishes correctly?
- Is there ever going to be a correct answer? (most likely no)

No 'perfect' pipeline – learn through experience



If unsure – always check **benchmark** studies

- Don't run programs that you are not sure the concepts
- Programs need to be **benchmarked**
- **Always look for most recent (and fair) benchmarks**

Bradnam *et al.* *GigaScience* 2013, **2**:10
<http://www.gigasciencejournal.com/content/2/1/10>



RESEARCH

Open Access

Assemblathon 2: evaluating *de novo* methods of genome assembly in three vertebrate species

Resource

Assemblathon 1: A competitive assessment of *de novo* short read assembly methods

Dent Earl,^{1,2} Keith Bradnam,³ John St. John,^{1,2} Aaron Darling,³ Dawei Lin,^{3,4} Joseph Fass,^{3,4} Hung On Ken Yu,³ Vince Buffalo,^{3,4} Daniel R. Zerbino,² Mark Diekhans,^{1,2} Ngan Nguyen,^{1,2} Pramila Nuwantha Ariyaratne,⁵ Wing-Kin Sung,^{5,6} Zemin Ning,⁷ Matthias Haimel,⁸ Jared T. Simpson,⁷ Nuno A. Fonseca,⁹ İnanç Birol,¹⁰ ...

Problem of not familiar with theories behind your program



Xkcd

Two situations you are most likely to encounter

Genome reference is available (for example, humans):

- Re-sequence (DNA, RNA)
- **Map** sequence to the genome

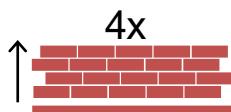
Genome reference is NOT available

- **Assemble** the reads to get the genome

New situation recently arise:

- Assemble the reads to produce many genomes of the same species

More Definition

 50-500 bp	Read	A sequenced piece of DNA
300-600 bp insert 	Paired-end read	Sequencing both ends of a short DNA fragment
> 1 kbp insert 	Mate-pair read	Sequencing both ends of a long DNA fragment
 length	Insert size	The length of the DNA fragment
	Contig	A set of overlapping DNA segments that represents a consensus region of DNA
 N 	Scaffold	Contigs separated by gaps of known length
 4x	Coverage	The number of times a specific position in the genome is covered by reads

What is an alignment?

Align the following two sequences:

ATTGAAAGCTA
GAAATGAAAAGG
1:
--ATTGAAA-GCTA
| | | | |
GAAATGAAAAGG--

2:
ATTGAAA-GCTA---
| | | | |
---GAAATGAAAAGG

Scoring scheme is needed:
1 for match
-1 for mismatch
-2 for gap

insertions / deletions (indels) mismatches
Which alignment is better?

Assembly (Lecture 4)



Genome
(3.000.000 letters)

Sequencing



Reads
(50-500 letters each)

Assembly



Genome
(3.000.000 letters)

A bad assembly (Lecture 4)



Genome
(3.000.000 letters)

Sequencing



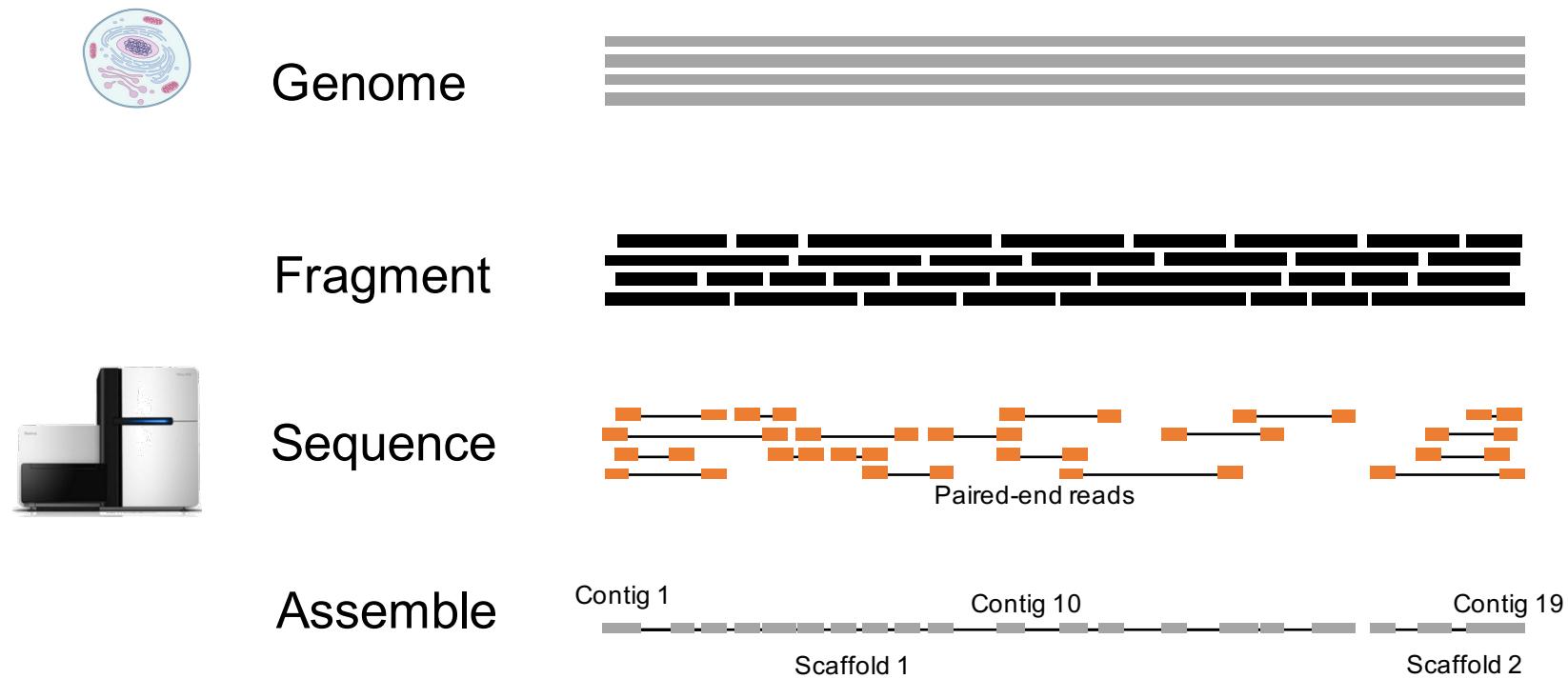
Reads
(50-500 letters each)

Assembly

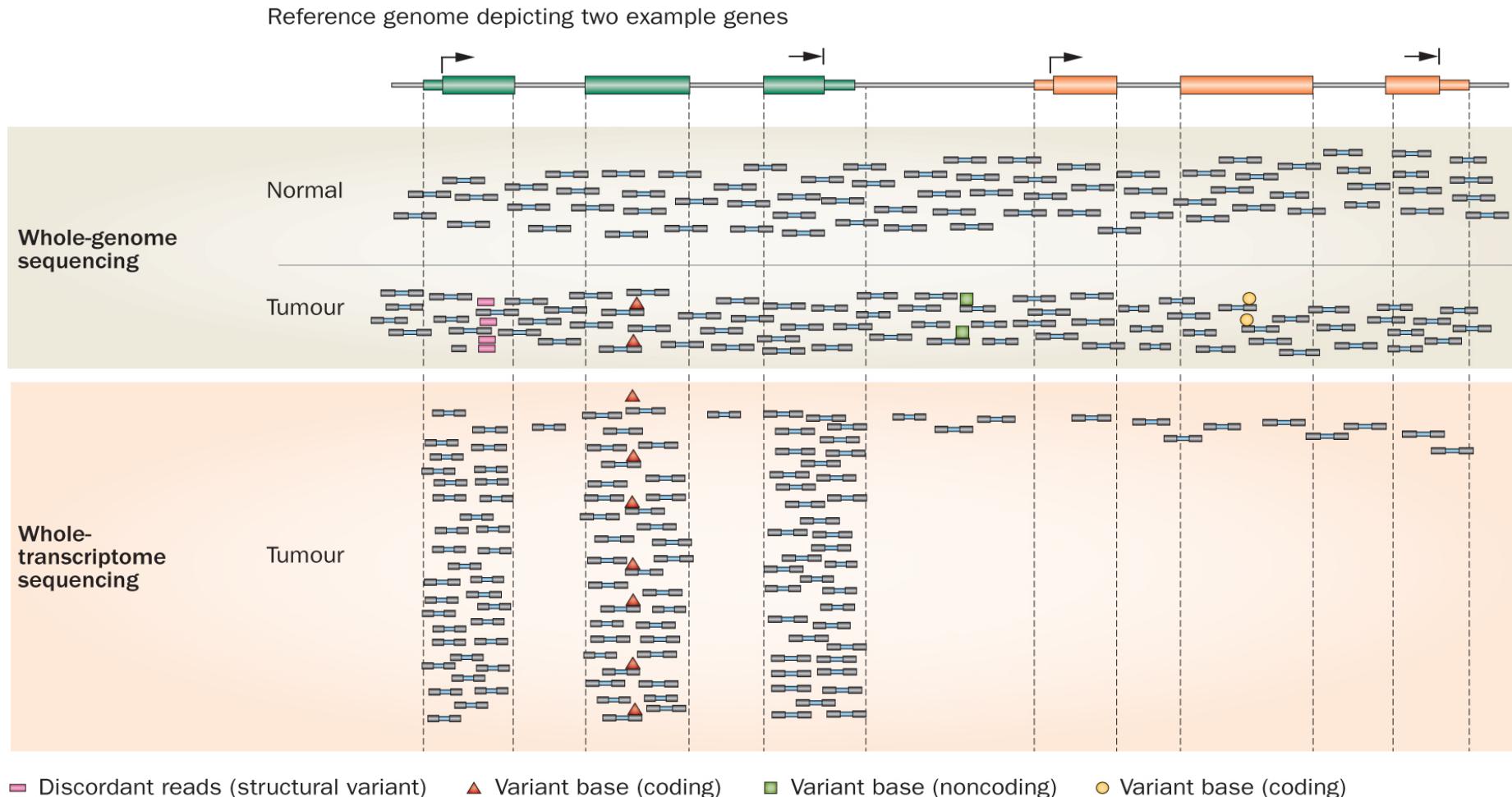


Genome
(3.000.000 letters)

Assembly (Lecture 4)



Mapping (Lecture 3)



doi:10.1038/nrgastro.2012.126

Read length matters in sequencing

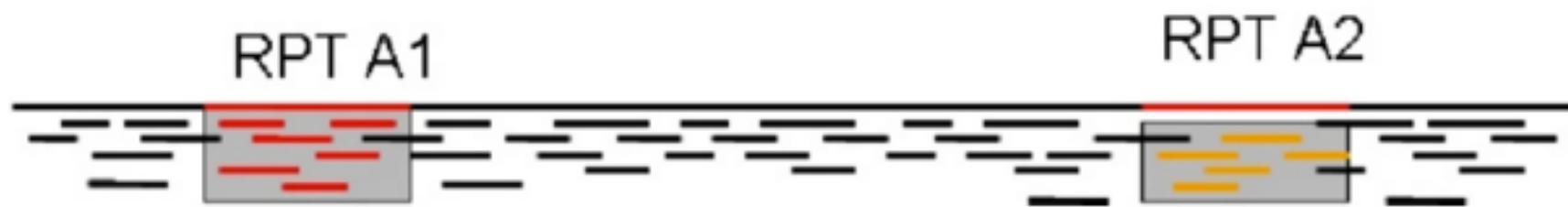


Figure 5. Two copies of a repeat along a genome. The reads colored in red and those colored in yellow appear identical to the assembly program.

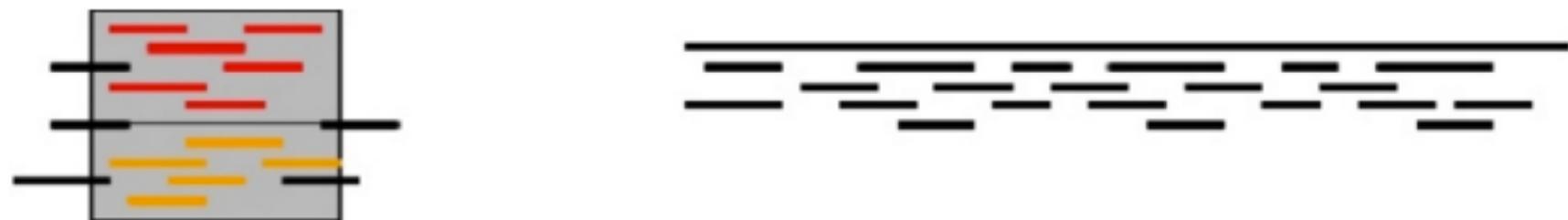
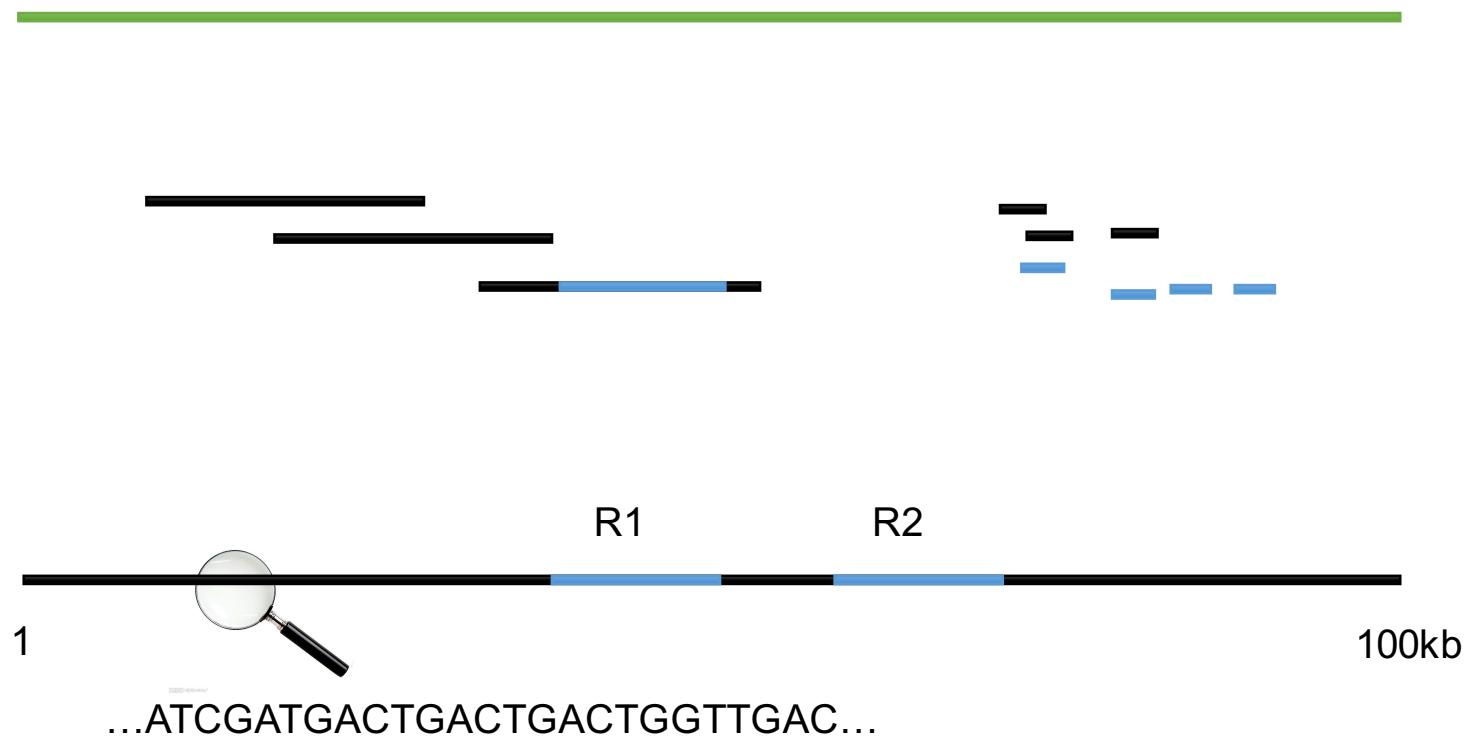
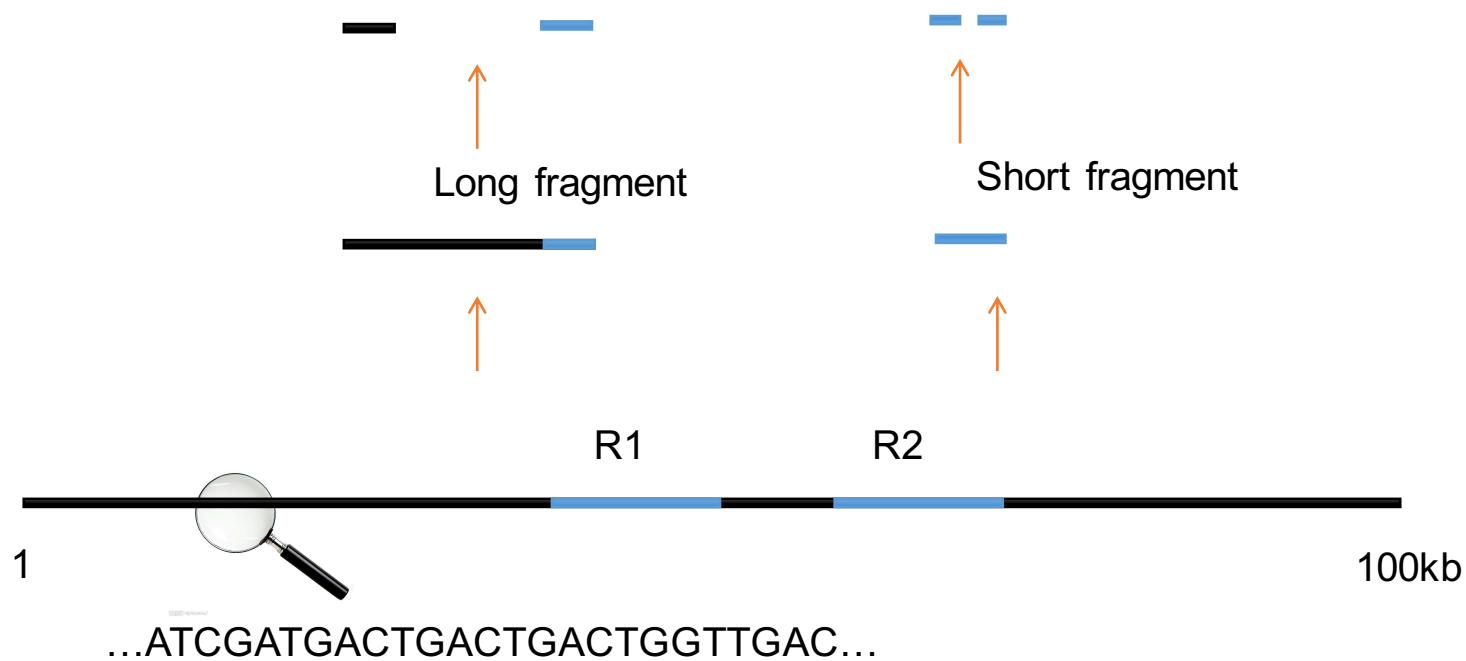


Figure 6. Genome mis-assembled due to a repeat. The assembly program incorrectly combined the reads from the two copies of the repeat leading to the creation of two separate contigs

Read length matters in sequencing



Paired end and insert size matter in sequencing



Depth matters in sequencing

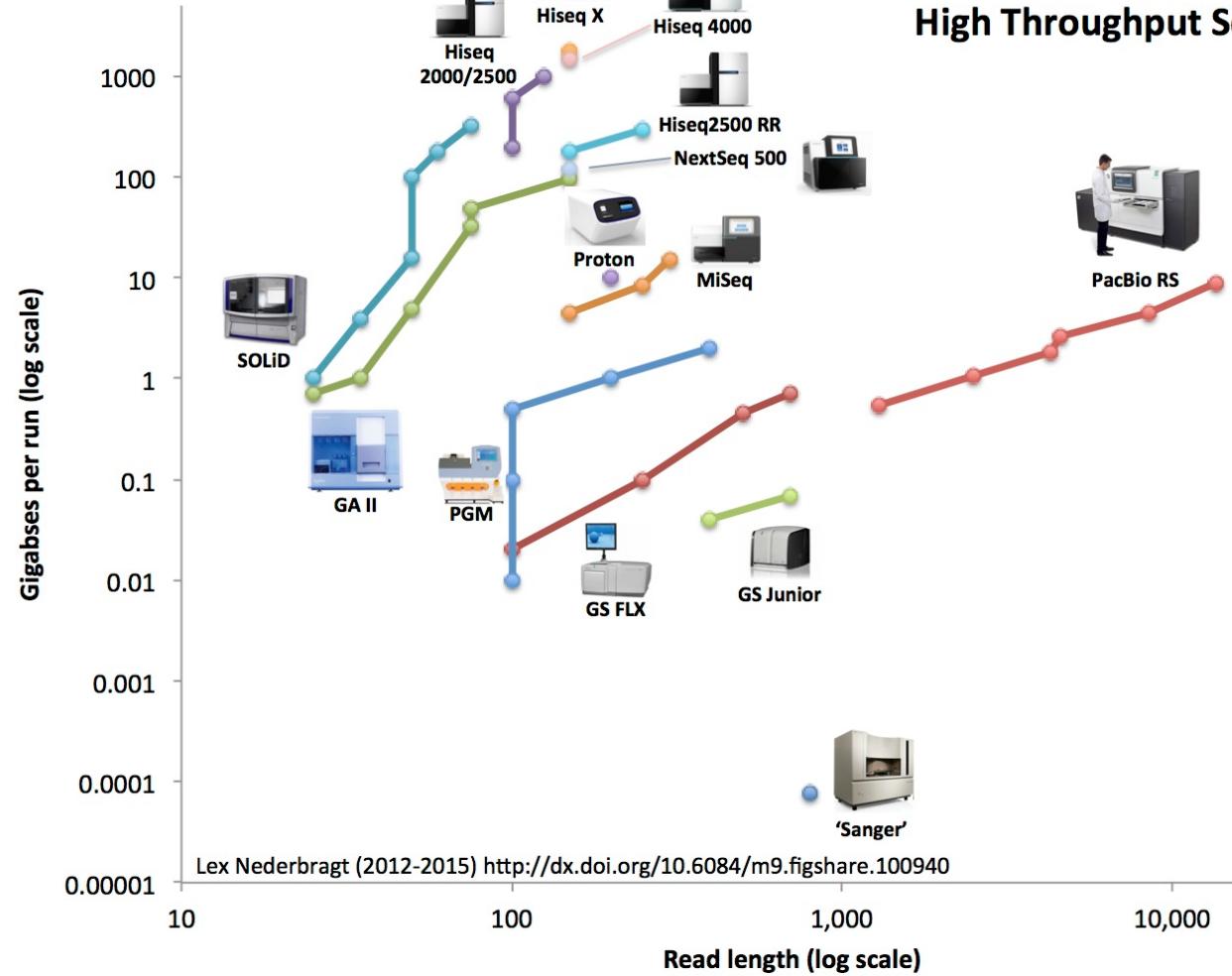
ATCGATGACTGACTGA~~A~~TGGTTGAC
ATCGATGACTGACTGA~~A~~TGGTTGAC
ATCCATGACTGACTGA~~A~~TGGTTGAC
ATCGATGACTGACTGA~~A~~TGGTTGAC
ATCGATGACTGACTGA~~A~~TGGTTGAC
ATCGATGACTGAGTGA~~A~~TGGTTGAC
ATCGATGACTGAGTGA~~A~~TGGTTGAC
ATCGATGACTGAGTGA~~A~~TGGTTGAC
ATCGATGACTGAGTGA~~A~~TGGTTGAC
10X ATCGATGACTGAGTGA~~A~~TGGTTGAC

1X ATCGAT~~C~~ACTGACTGACTGGTTGAC
Homozygous? Heterozygous?

...ATCGATGACTGACTGACTGGTTGAC...

reference

Developments in High Throughput Sequencing



Different sequencing platforms

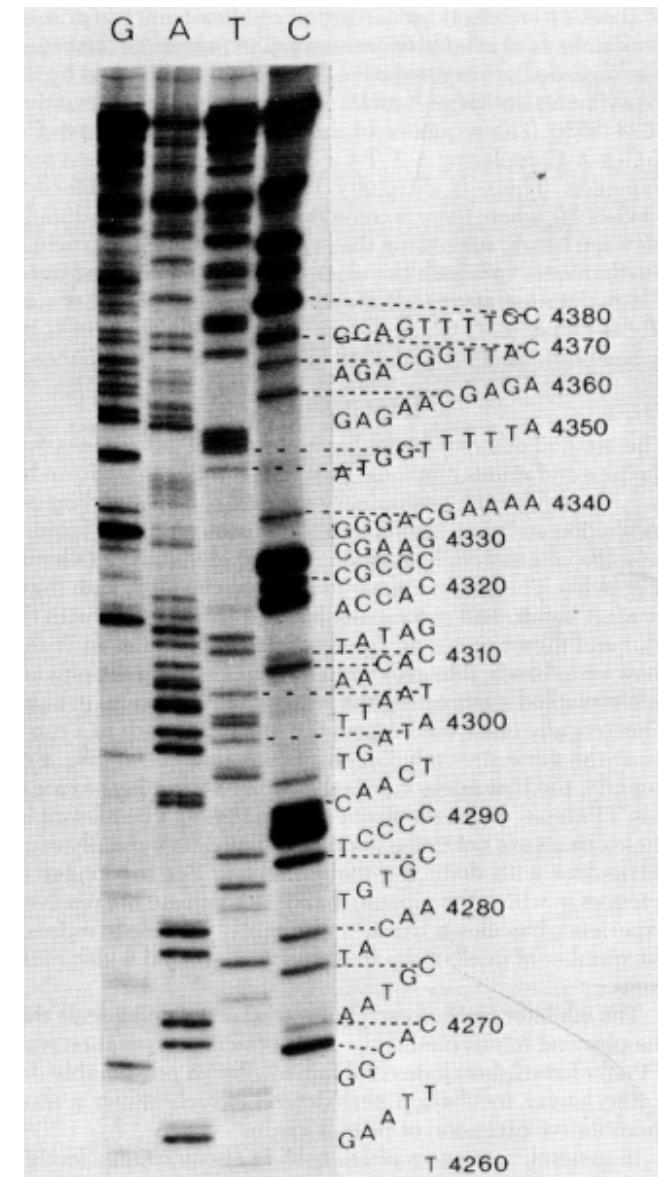
DNA sequencing with chain-terminating inhibitors

(DNA polymerase/nucleotide sequences/bacteriophage ϕ X174)

F. SANGER, S. NICKLEN, AND A. R. COULSON

Medical Research Council Laboratory of Molecular Biology, Cambridge CB2 2QH, England

Contributed by F. Sanger, October 3, 1977

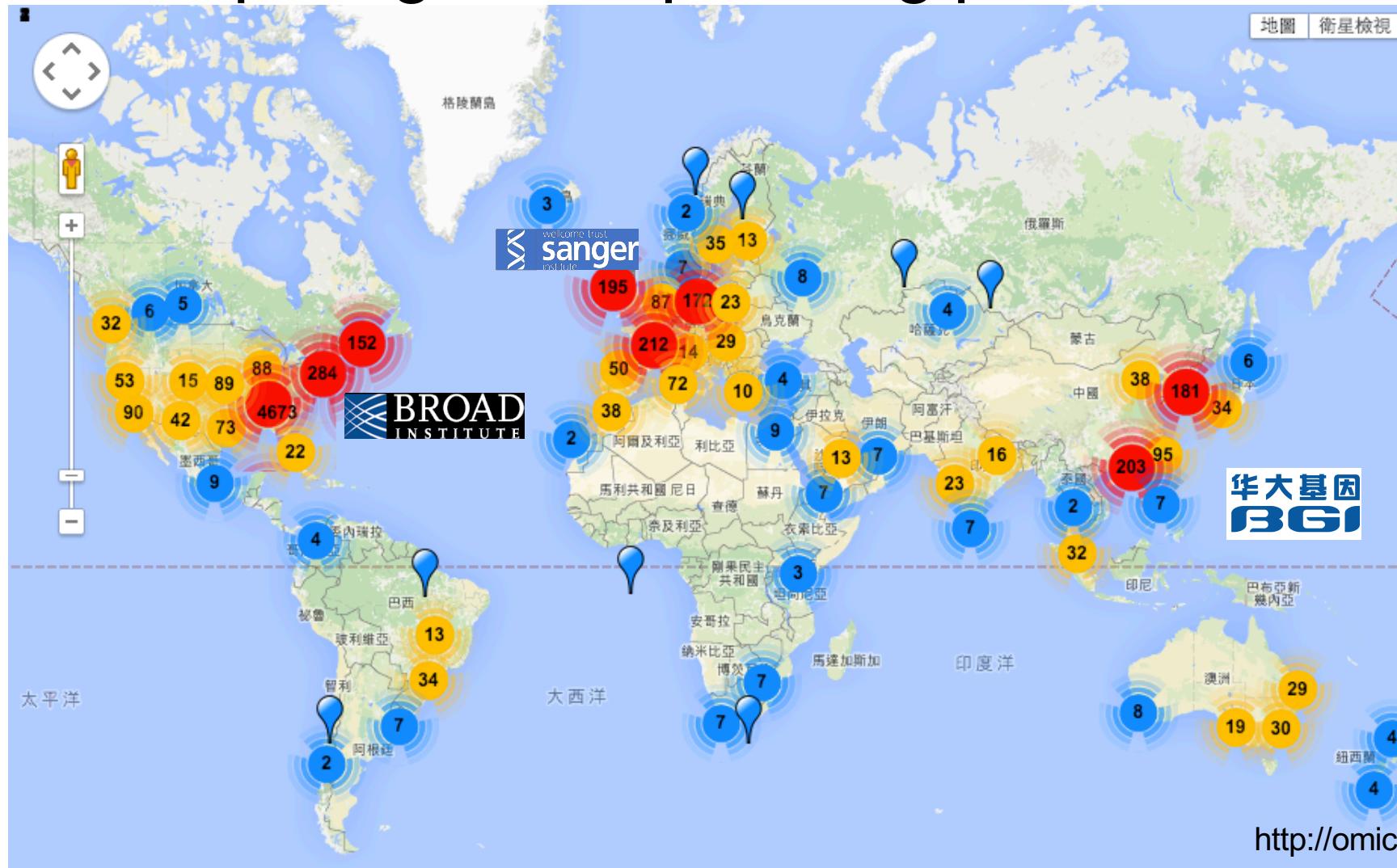


ABI 3730xi at TIGR



<https://www.flickr.com/photos/jurvetson/57080968>

World competing for sequencing power

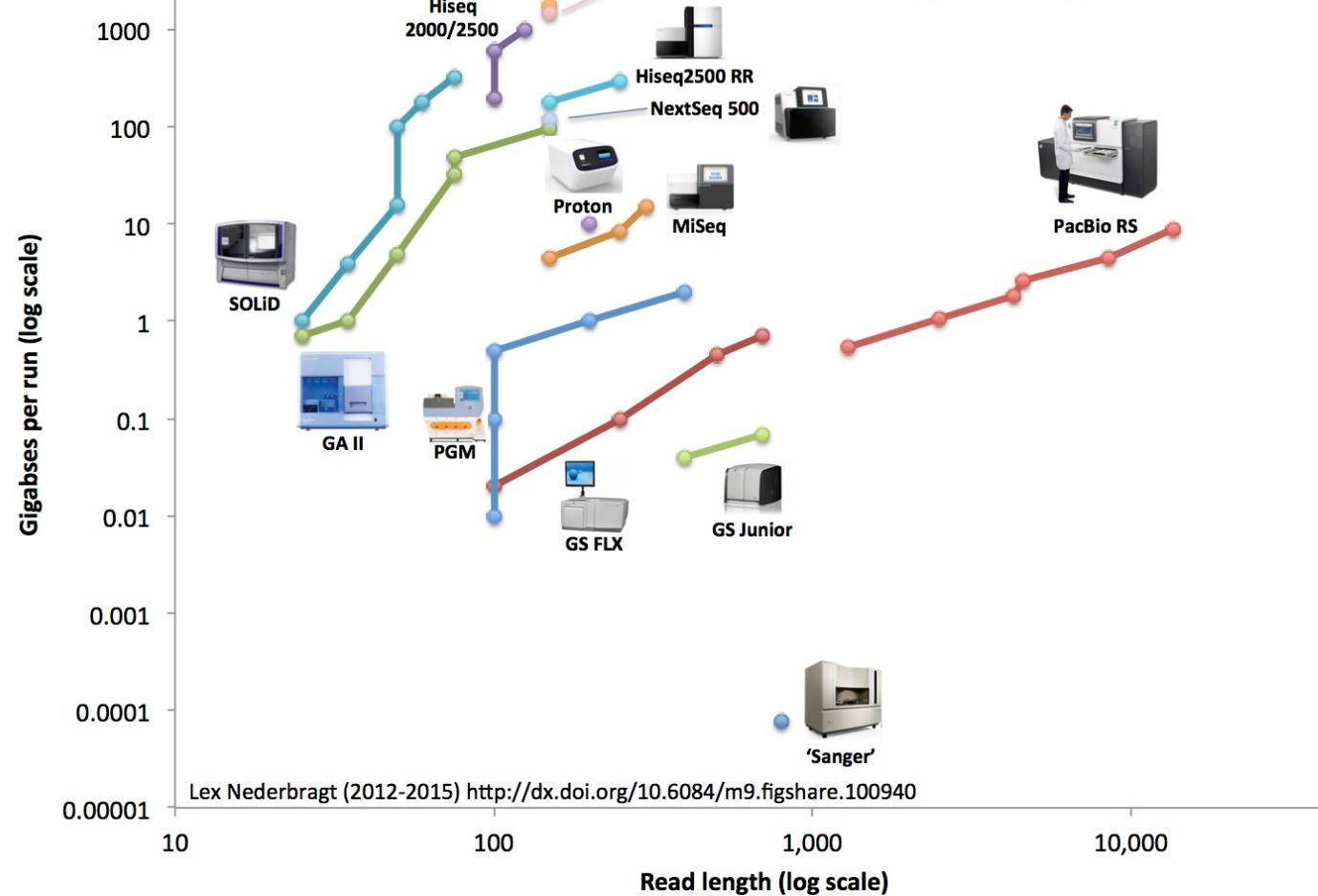


<http://omicsmaps.com/>

Sequencing Platforms

- Short reads
 - 1. ~~Genome Analyzer IIx (GAIIx)~~ – Illumina
 - 2. HiSeq2000, HiSeq2500, MiSeq – Illumina
- Long reads
 - 1. ~~Genome Sequencer FLX System (454)~~ – Roche
 - 2. PacBio RS - Pacific Bioscience
 - 3. GridION – Oxford Nanopore

Developments in High Throughput Sequencing



Platform	GS jnr	FLX plus	MiSeq	Next Seq 500	HiSeq 2500 RR	Hiseq 2500 V3	HiSeq 2500 V4	HiSeq 4000	HiSeq X	SOLiD 4	5000 XL	318 HiQ 520	Ion 530	Proton P1	PGM HiQ 540	RS P6-C4	Sequel	Mini ION	Prome thION	QiaGen Gene Readr	BGI SEQ 500	#	
Reads: (M)	0.1	1.25	25	400	600	3000	4000	5000	6000	1400	--	5	15-20	165	60-80	5.5	38.5	0.05	--	400	--	--	
Read length: (paired-end*)	400	700	300*	150*	100*	100*	125*	150*	150*	50	75	200 400	200 400	200	220	15K	12K	10K	10K	--	--	--	
Run time: (d)	0.4	0.9	2	1.2	1.125	11	6	3.5	3	12	7	0.37	--	--	--	4.3	4.3	2	--	--	1	--	
Yield: (Gb)	0.035	0.7	15	120	120	600	1000	1500	1800	100	180	1.2-2	6-8	10	10-15	12	84	0.5	600	80	200	--	
Rate: (Gb/d)	0.2	0.75	7.5	100	106.6	55	166	400	600	8.3	30	--	--	--	--	2.8	19.5	0.25	--	--	--	--	
Reagents: (\$K)	1.1	6.2	1	4	6.145	23.47	29.9	29.9	12.75	9	10.5	0.6	--	1	1.2	2.4	11.2	1	--	0.5	--	--	--
per-Gb: (\$)	31K	8K	93	33.3	51.2	39.1	29.9	20	7	90	58.33	--	--	100	--	200	80	2000	20	--	--	--	--
Machine: (\$)	110K	500K	99K	250K	740K	690K	690K	900K	1M	500K	595K	50K	65K	243K	242K	695K	350K	1000	30K	--	--	--	--

#Page maintained by <http://www.vilellagenomics.com> #Editable version: tinyurl.com/ngsspecsshared

```
#curl "https://docs.google.com/spreadsheets/d/1GMMfhYLK0-q8Xklo3YxiWaZA5vVMuhU1kg41g4xLkXc/export?gid=4&format=csv" | grep -v '^#' | column -t -s\, | less -S
```



Albert Vilella @AlbertVilella · 18h

Updated NGS specs @BGI_Genomics @PacBio @illumina @thermofisher @nanopore @QIAGENscience tinyurl.com/ngsspecs



1

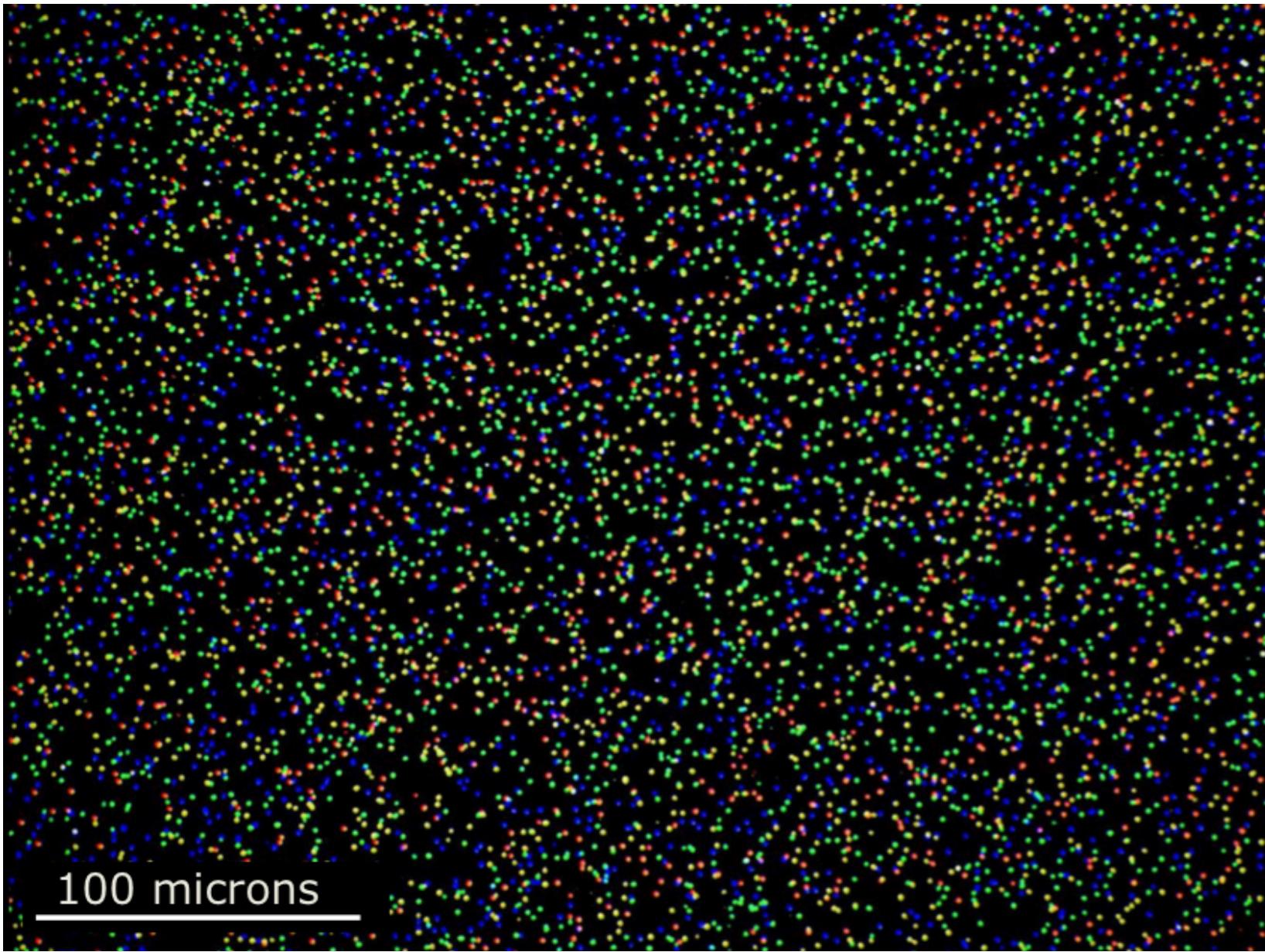


5



Illumina HiSeq





HiSeq and MiSeq

HiSeq 2000

Initially capable of up to 600Gb per run in 13 days.

Cost of resequencing one human genome:

30x coverage about \$6,000- \$9,000



HiSeq 2500

Initially capable of up 100Gb per run in 27hours.

Cost per genome - ???

MiSeq

-Small capacity system. PE 2x250cycles in 24 hours.

-Long insert size possible: 1.5kb, 3kb

-2x400bp in R&D



HiSeq X Ten!

Population Power.



HiSeq X_{TEN}

Ultra-high-throughput.
Population scale projects.

The next revolution in sequencing has arrived. HiSeq X Ten is the first sequencing platform to break the \$1000 barrier for a 30x human genome. HiSeq X Ten is a set of ten ultra-high-throughput sequencers, built to sequence tens of thousands of human whole genomes per year.

The HiSeq X Ten contains 10 sequencing systems.

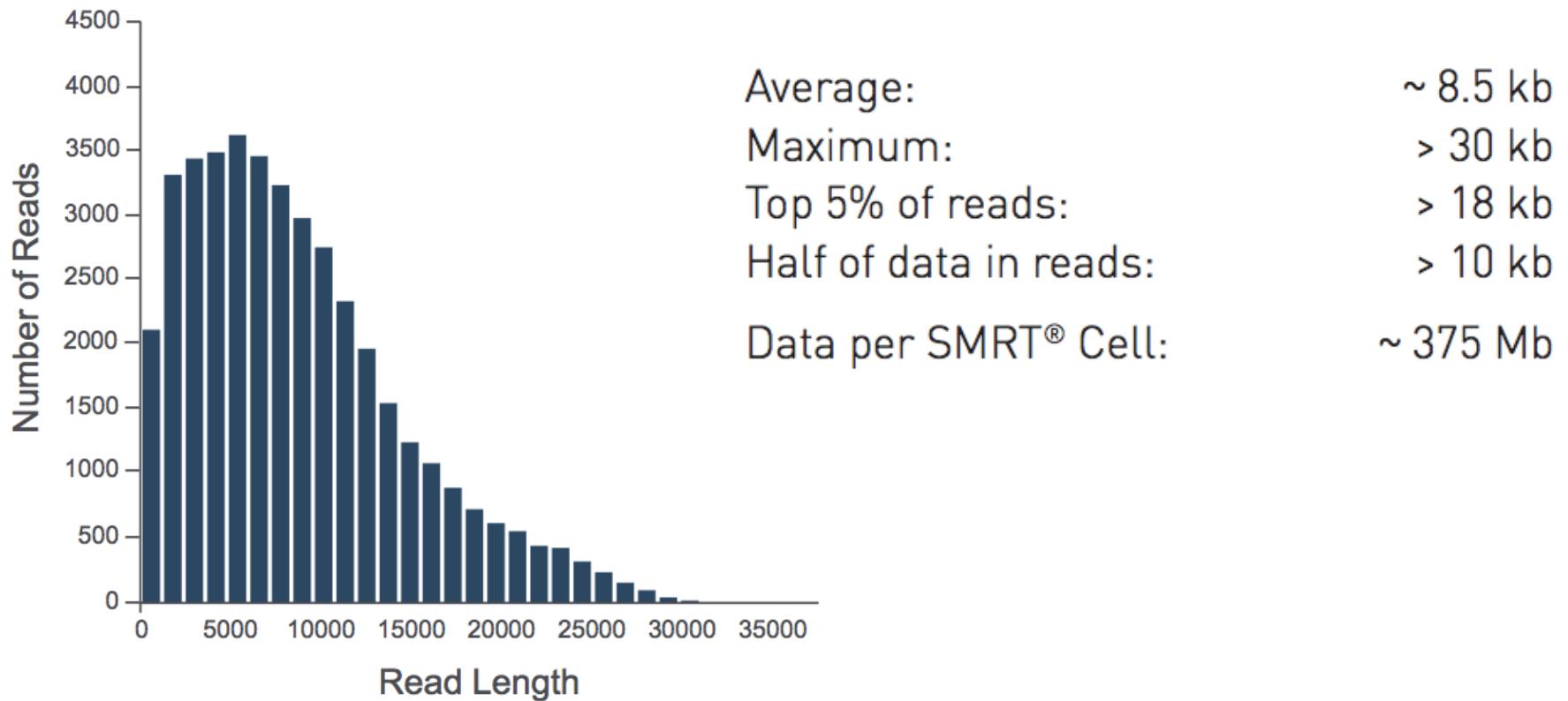
<http://www.illumina.com/systems/hiseq-x-sequencing-system.ilnn>

PacBio (Pacific Biosciences)

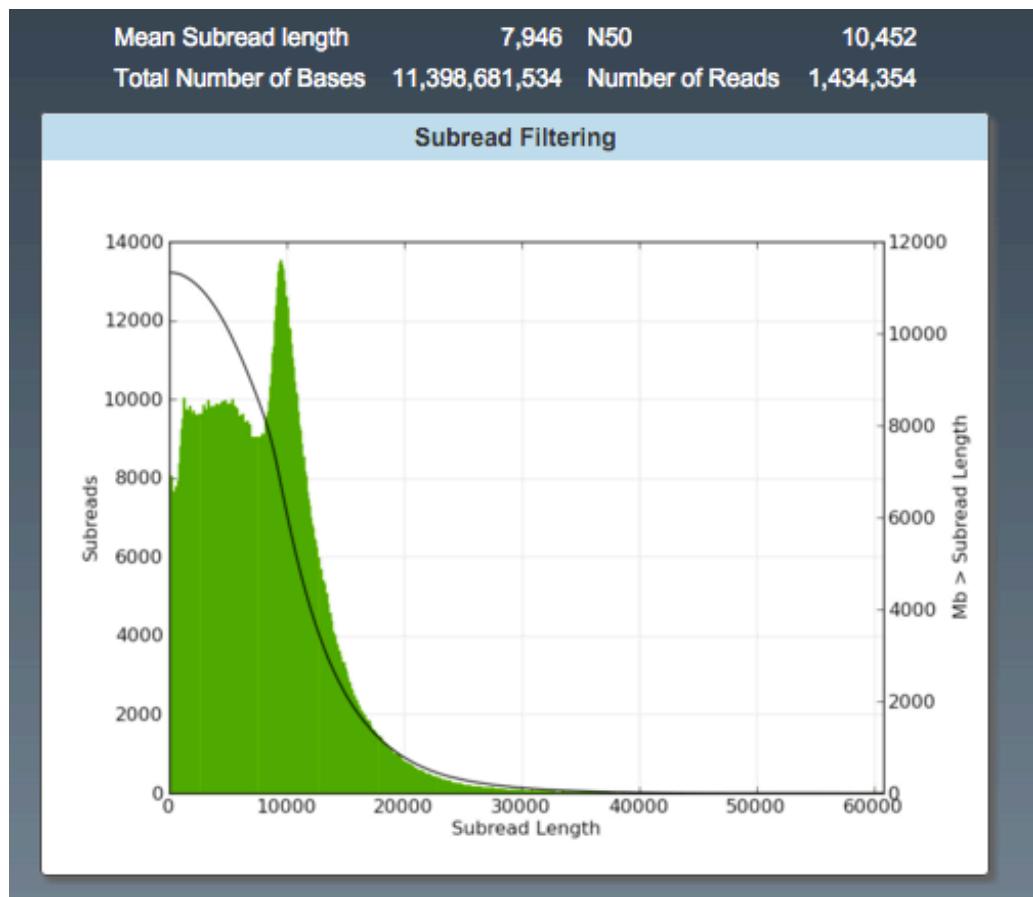
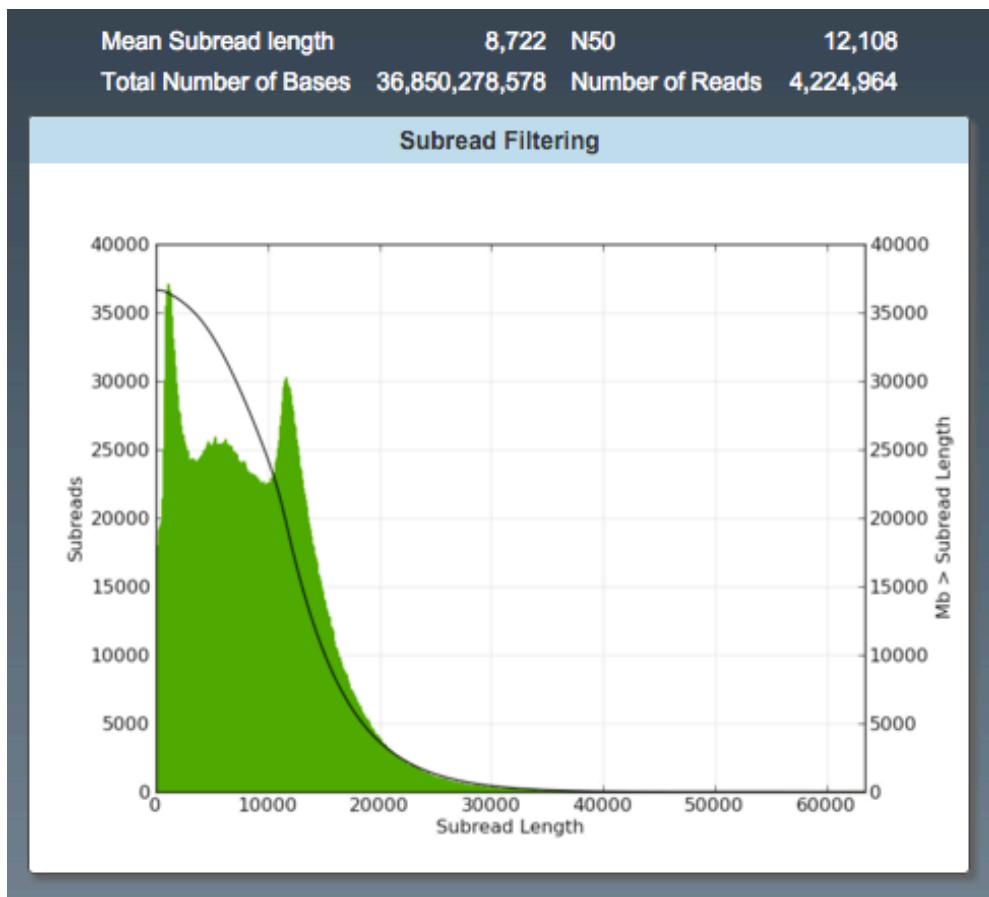


http://files.pacb.com/pdf/PacBio_RS_II_Brochure.pdf

PacBio (Pacific Biosciences)



PacBio (Pacific Biosciences)

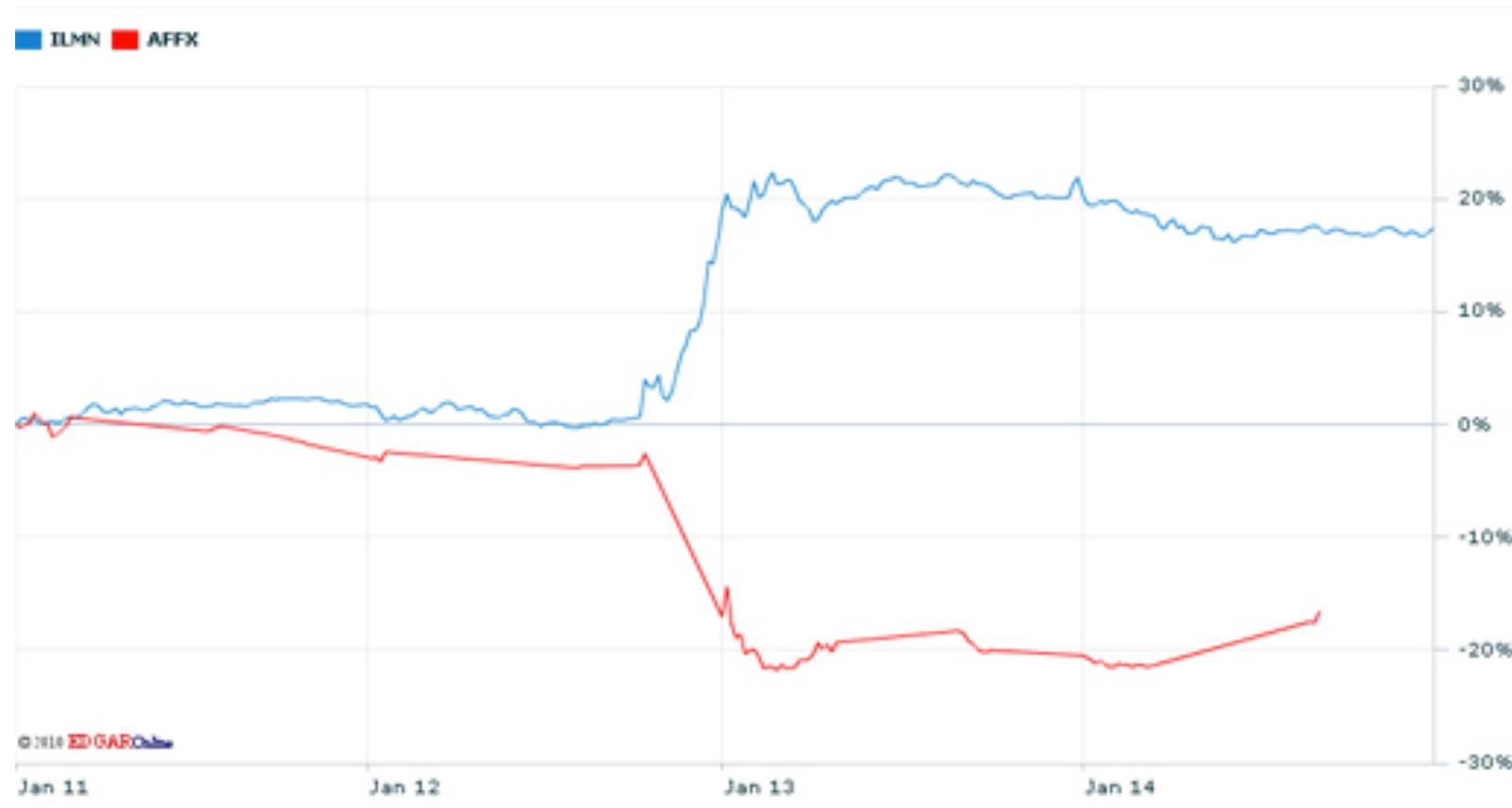


Oxford Nanopore



- 150Mb per run
- Tested 48kb read length
- \$900 per instrument
- 500 pores per device
- XXXMb per run
- Tested 48kb read length
- \$XXX per instrument
- 2000 pores per device, soon 8000 pores
- Cost per human genome \$1500.

Come and go of technologies



Data types

A lot of data

- We biologists generate a lot of data
 - Experiments, sequencing
 - Everything is more high throughput, but not necessarily less noisy
- Different data types
 - Images, Sequences, Signals, Locations, Linkage, Frequencies...
- How do we
 - analyse them?
 - store them?
 - publish them?
 - reuse them?

Always understand your data / programs

- Understand:
 - Data format
 - The nature of your data
- Please don't
 - assume data you are given is 'correct'
 - Scenario 1: We got the assmeblies and analysis from company XXXX, and we don't know what to do with it
 - assume everything's correct online
 - Run everything in 'default' mode

FASTA format

```
>Name_of_sequence
GCAGGGCATCCGCTGCGTGCTGGGCAAAGTCTGTGGCGA
CCGTATGAAACCCCTGAACCCGGCAATACGGTTGCAGATG
CTGGAGGCAGGCGGCAACGGCAATCAGCTTGATTGAGGT
GAGGTACATAAGCCCCGCCAGGCGAGCAGGCAAAGCTGT
TCTGGTCCGAGCCCCAACAGGGTTCACCAAGTGGCCTG
CCGACGAAAGCGCCGAAGCCCCAACCCATCAATCACTG
GCAGCGTGCAGTCCAGGCCATCGACGAGGCCATCATTGA
AGCGCGGTACGACCCCCGAAACGGCACGCTCATTGTTGC
GTTGGCTTCCTATGGTCGGCGCGACCCAGCTTCCCTGGA
ACAGTTGCGGCCACCTCGCGAAGGAAGGCATTCCCC
CGGAATTCTGTACACATTATGAGCCTGACGGACCCTTGC
```

Alignment format

- Some programs need slightly modified format

```
>Name_of_sequence_1
GCAGGGCATCCGCTGCGTGTGGCAAAGTCTGTGGCGA
CCGTATGAAACCCCTGAACCCGGCAATACGGTTGCAGATG
CTGGAGGCAGGCGGGCAACGGCAATCAGCTTGATTGAGGT
GAGGTACATAAGCCCCGCCAGGCGAGCAGGCAAAGCTGT
TCTGGTCCGAGCCCCAACAGGGTTCACCAAGTGGCCTG
CCGACGAAAGCGCCGAAGCCCG
>Name_of_sequence_2
GCAGGGCATCCGCTGCGTGTGGCAAAGTCTGTGGCGA
CCGTATGAAACCCCTGAACCCGGCAATACGGTTGCAGATG
CTGGAGGCAGGCGGGCAACGCAATCAGCTTGATTGAGGTG
AGGTACATAAGCCCCGCCAGGCGAGCAGGCAAAGCTGTTC
TGGTCCGAGCCCCAACAGGGTTCACCAAGTGGCCTGCC
GACGAAAGCGCCGAAGCCCG
```

Data type keep evolving

- Very first fastq file was invented in 2007?
- Obviously will become problematic in storage later on...

>Name_of_sequence_1

GCGGGTA

>Name_of_sequence_1

20 30 33 30 20 33 19

Fastq files:

FASTQ format is a **text-based format** for storing both a biological **sequence** (usually nucleotide sequence) and its corresponding **quality scores**.

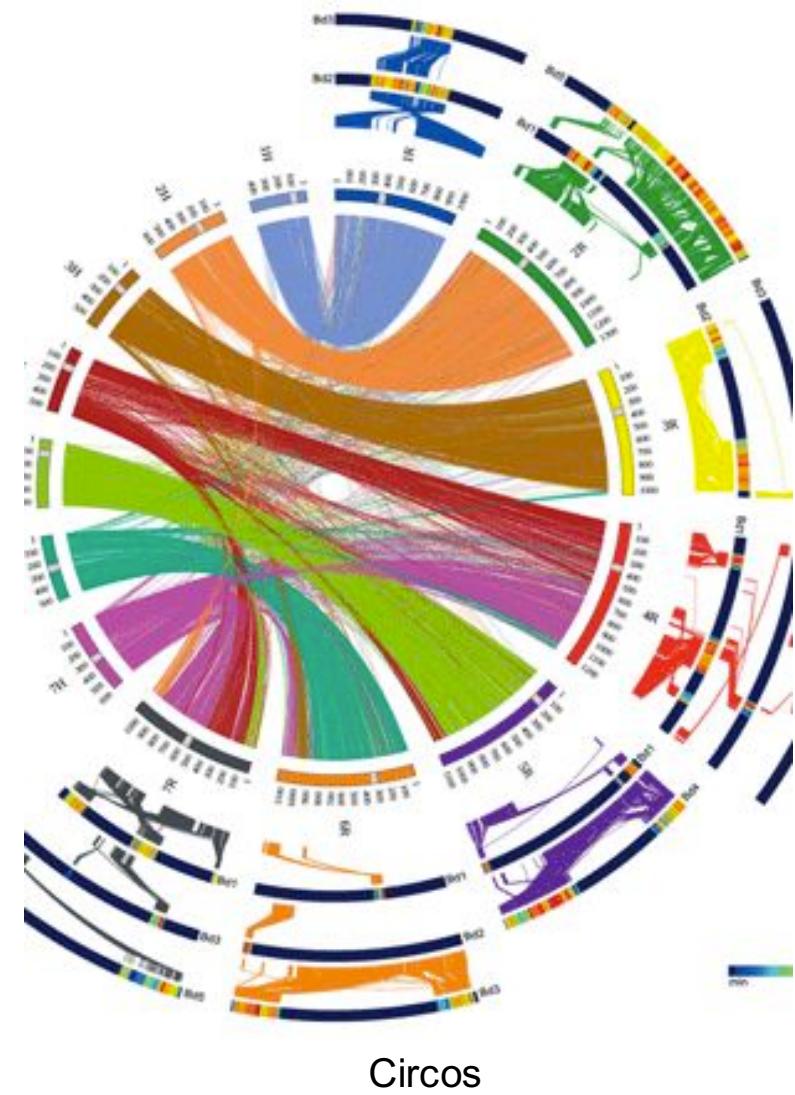
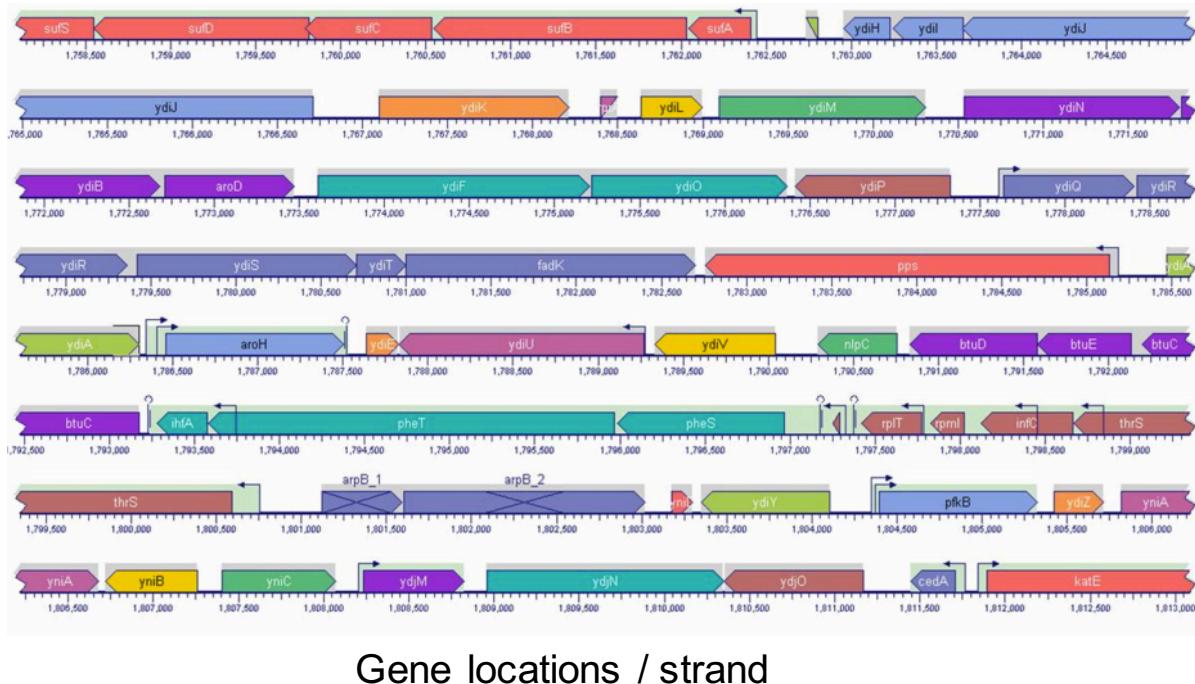
-Wikipedia

```
@SEQUENCE_ID1
ATGCGCGCGCGCGCGCGCGGGTAGCAGATGACGACACAGAGCGAGGATGCGCTGAGAGTA
GTGTGACGACGATGACGGAAAATCAGA
+
BBBBBPPPPPXXXXX^~~~~~^~~~~~^~~~~~^~~~~~^_ ~~~~~^~~~~~^_eeeeeee
[[[[[ ^^^ ] ] ]XXXXXPPPPPBBB
```

1. Single line ID with at symbol (“@”) in the first column.
2. There should be not space between “@” symbol and the first letter of the identifier.
3. Sequences are in multiple lines after the ID line
4. Single line with plus symbol (“+”) in the first column to represent the quality line.
5. Quality ID line can have or have not ID
6. Quality values are in multiple lines after the + line

Locations / maps

- How do we represent/visualise them?



BED/gff format

- Features on genome use bed / gff files to represent their locations
- “Optional field” can be added for additional information

```
chr7 127471196 127472363
chr7 127472363 127473530
chr7 127473530 127474697
chr7 127474697 127475864
chr7 127475864 127477031
chr7 127477031 127478198
chr7 127478198 127479365
chr7 127479365 127480532
chr7 127480532 127481699
```

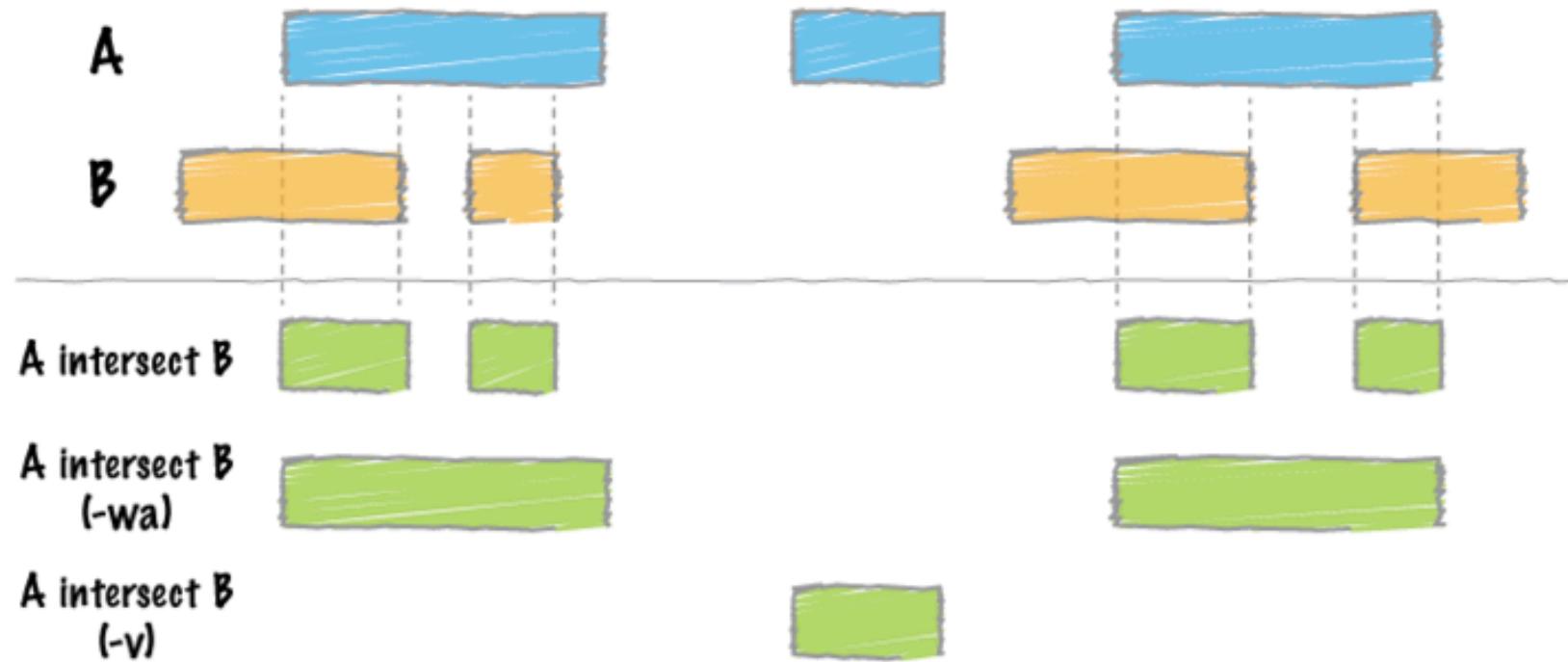
IV	curated exon	5506900	5506996	.	+	.	Transcript	B0273.1
IV	curated exon	5506026	5506382	.	+	.	Transcript	B0273.1
IV	curated exon	5506558	5506660	.	+	.	Transcript	B0273.1
IV	curated exon	5506738	5506852	.	+	.	Transcript	B0273.1

<http://genome.ucsc.edu/FAQ/FAQformat#format1>

<http://gmod.org/wiki/GFF2>

Bedtools – extremely useful

Intersect w/
1 database



SAM format

- 1 DNA is extracted from a sample.
- 2 DNA is sequenced.
- 3 Raw sequencing reads are aligned to a reference genome.
- 4 Aligned reads are evaluated and visualized.
- 5 Genomic variants, including single nucleotide polymorphisms (SNPs), small insertions and deletions are identified.

samtools

SAM format

- Everything in one line

```
HISEQ:134:C6H9FANXX:8:1110:10236:94013 99 chr1 11844 0 150M = 12057 363
GGTATCATTACCACTTTCTTCGTTAACCTGCCGTAGCCTTCTTGACCTCTTCTTGTTC
ATGTGTATTGCTGTCTTAGCCCAGACTCCCGTATCCTTCCACCAGGCCTTGAGAGGTCACA
GGGTCTTGATGCTG
>A=>AFDEEGEGGEFFFFFFFGCDFBEGFFHFGCDGEHGGFFFFGFFEDGGFGFFGFFFDFEDF
GCFHCHDBEFFHFEGCFEFED@CEEEBEADCBBCB>?,?AA@@@?@?>@?;?@??==?=?;<;?@GEH
GFDGFFHAC=?=@ MC:Z:150M
BD:Z:NNOOPSQQNOPMNOOMGGGNMMGGNONNLMNNONOPOMNPOPQOONHGONNHOPOOOO
NNONNHONPMNOPPPMNMONNHOPPPMQNONNM
```

<https://broadinstitute.github.io/picard/explain-flags.html>

<http://www.htslib.org/doc/samtools.html>

SAM format

- Bitwise flag

1	QNAME	Query template/pair NAME
2	FLAG	bitwise FLAG
3	RNAME	Reference sequence NAME
4	POS	1-based leftmost POSition/coordinate of clipped sequence
5	MAPQ	MAPping Quality (Phred-scaled)
6	CIGAR	extended CIGAR string
7	MRNM	Mate Reference sequence NaMe (=' if same as RNAME)
8	MPOS	1-based Mate POSition
9	TLEN	inferred Template LENgth (insert size)
10	SEQ	query SEQuence on the same strand as the reference
11	QUAL	query QUALity (ASCII-33 gives the Phred base quality)
12+	OPT	variable OPTIONAL fields in the format TAG:VTYPE:VALUE

mapped in correct orientation and within insert size

99	1+2+32+64	99	1	map	+	map	-	y
147	0+1+2+16+128	147	2	map	-	map	+	y

83	1+2+16+64	83	1	map	-	map	+	y
163	1+2+32+128	163	2	map	+	map	-	y

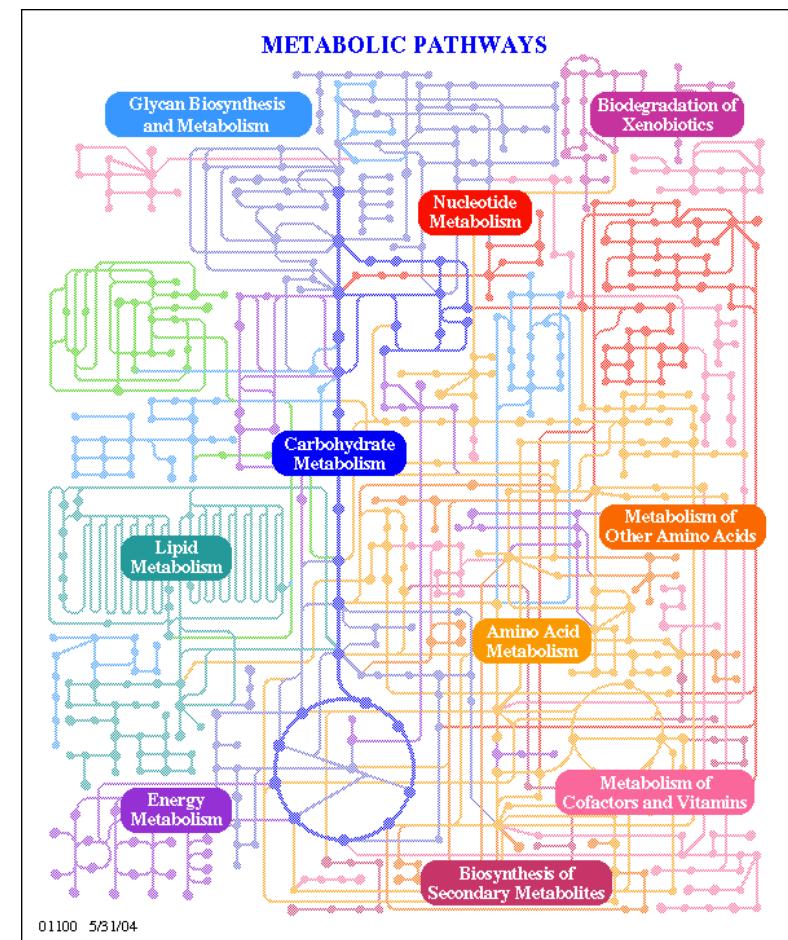
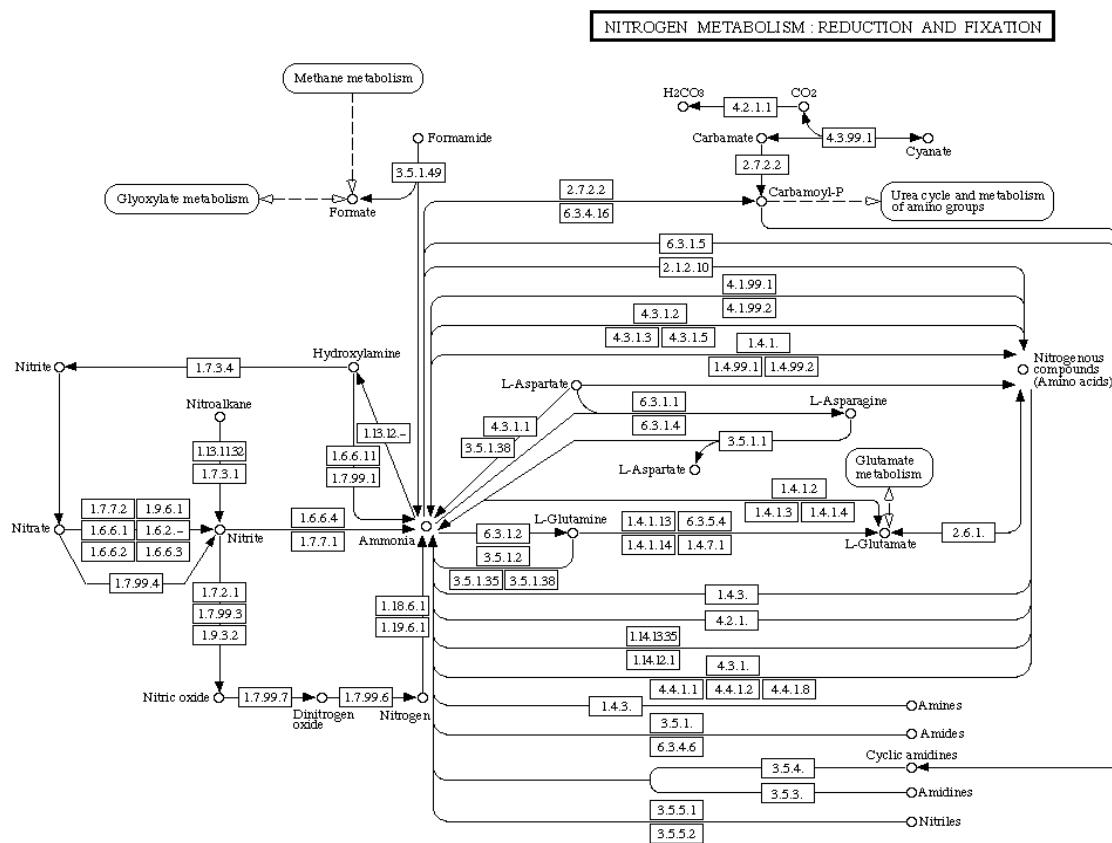
<https://broadinstitute.github.io/picard/explain-flags.html>

<http://www.htslib.org/doc/samtools.html>

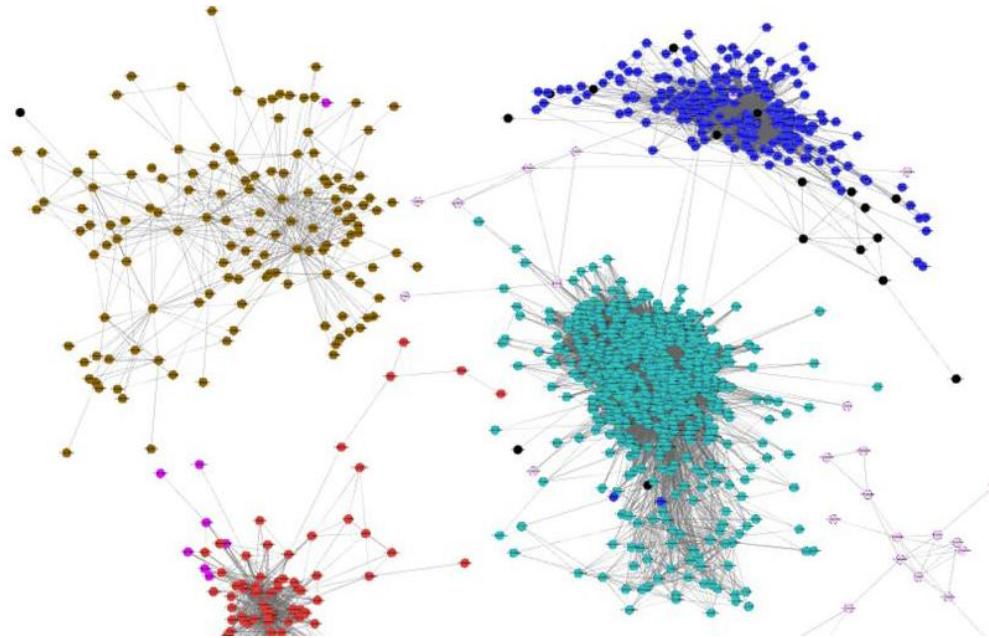
Public databases

- Access public data is key to bioinformatics analysis
 - We can't survive without pubmed
 - Any closely related species to your working species?
 - Any additional experimental data?
 - Any functional annotation to your SNP?
- Remember to deposit your own data to contribute

KEGG: Kyoto Encyclopedia of Genes and Genomes



Importance of networks in biology



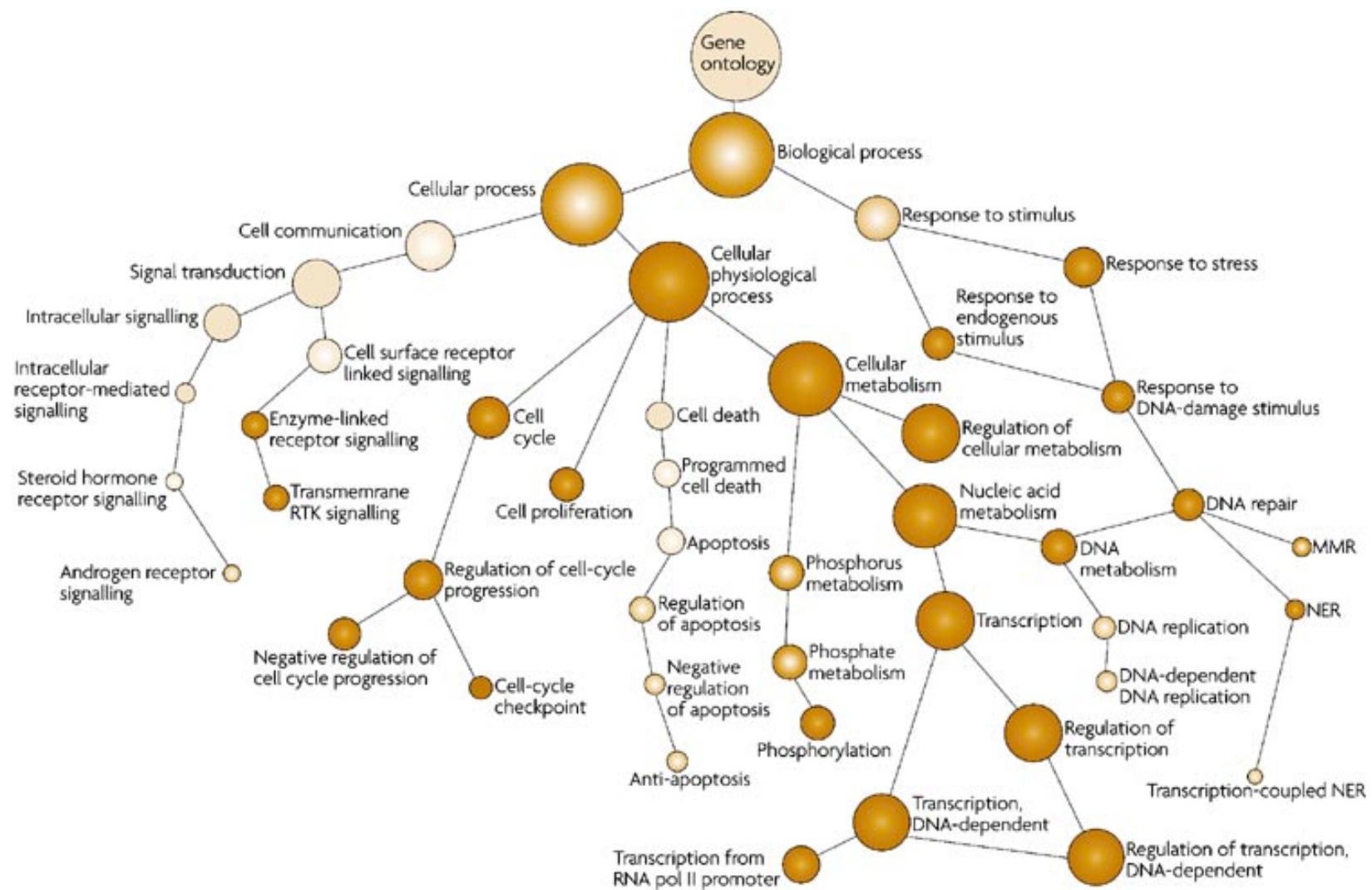
Gene interaction networks



Protein interaction network

Genome Ontology

- Biologists are fun:
 - “sonic hedgehog”
 - RING domain = Really Interesting New Gene
 - The ken and barbie gene
- An attempt to unify the names and functions across all species
- Genome Ontology uses a single 3 part system
 - Molecular function (specific tasks)
 - Biological process (broad biological goals - e.g. cell division)
 - Cellular component (location)



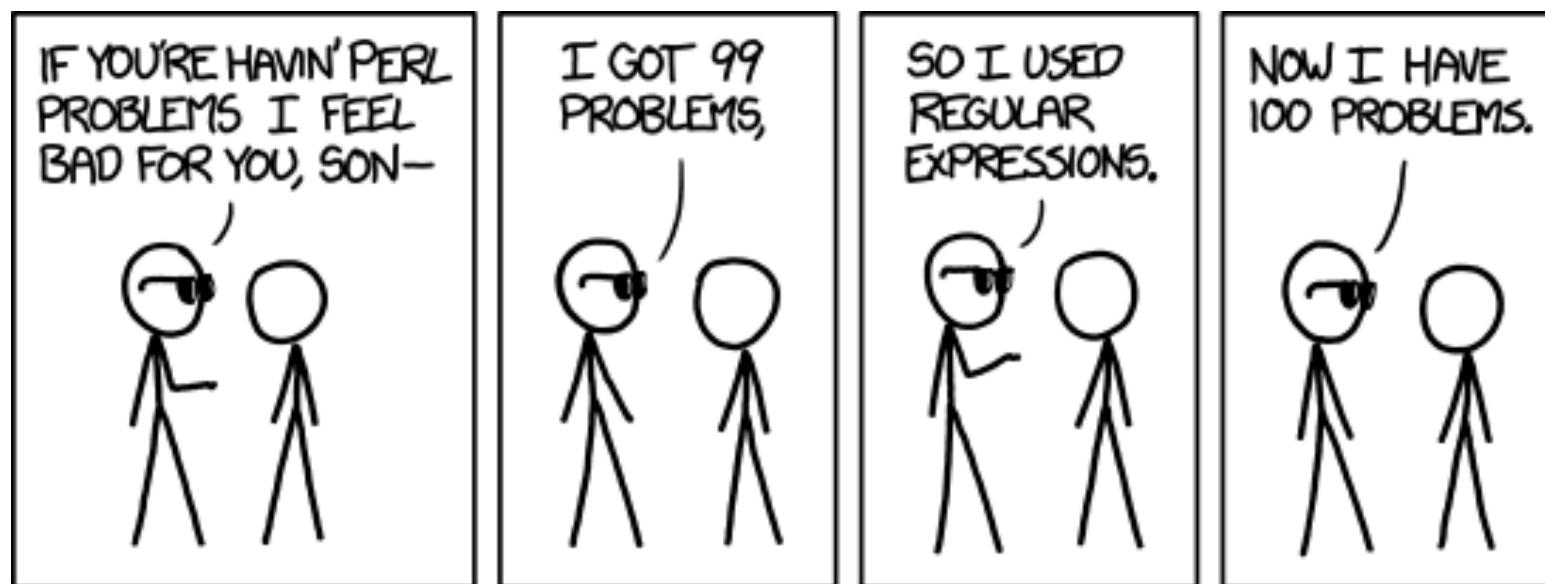
Most of these data are **text** based

- All output = text files
- So **familiarity** to deal with large amount of data is expected
- That's why **Perl** was popular back in the early period



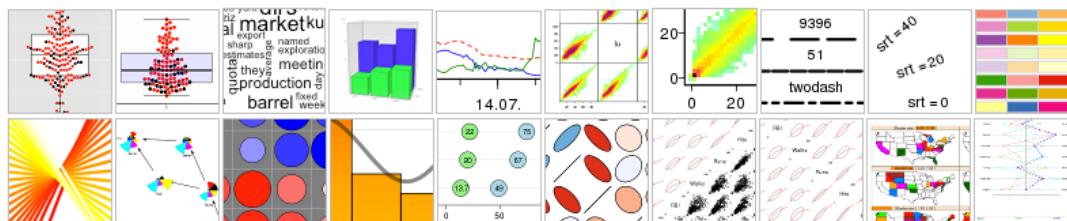
Problem with Perl

- Quick and ‘dirty’
- Data gets increasingly bigger and more complicated

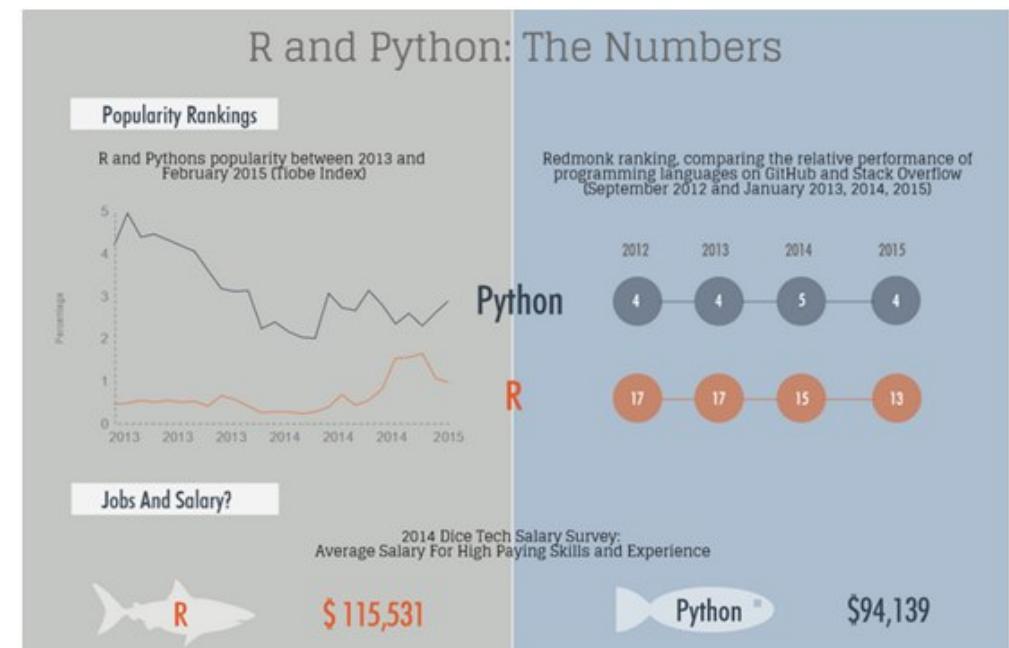
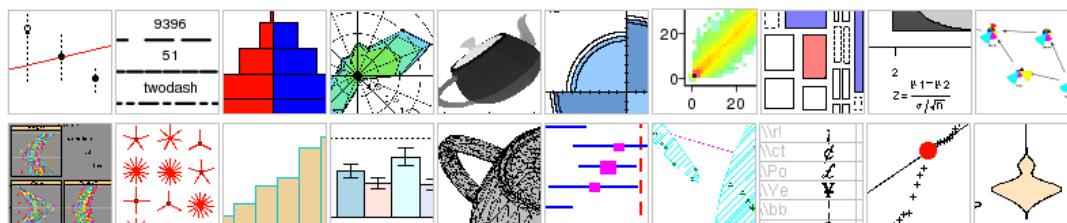


Python and R

» Last entries ...

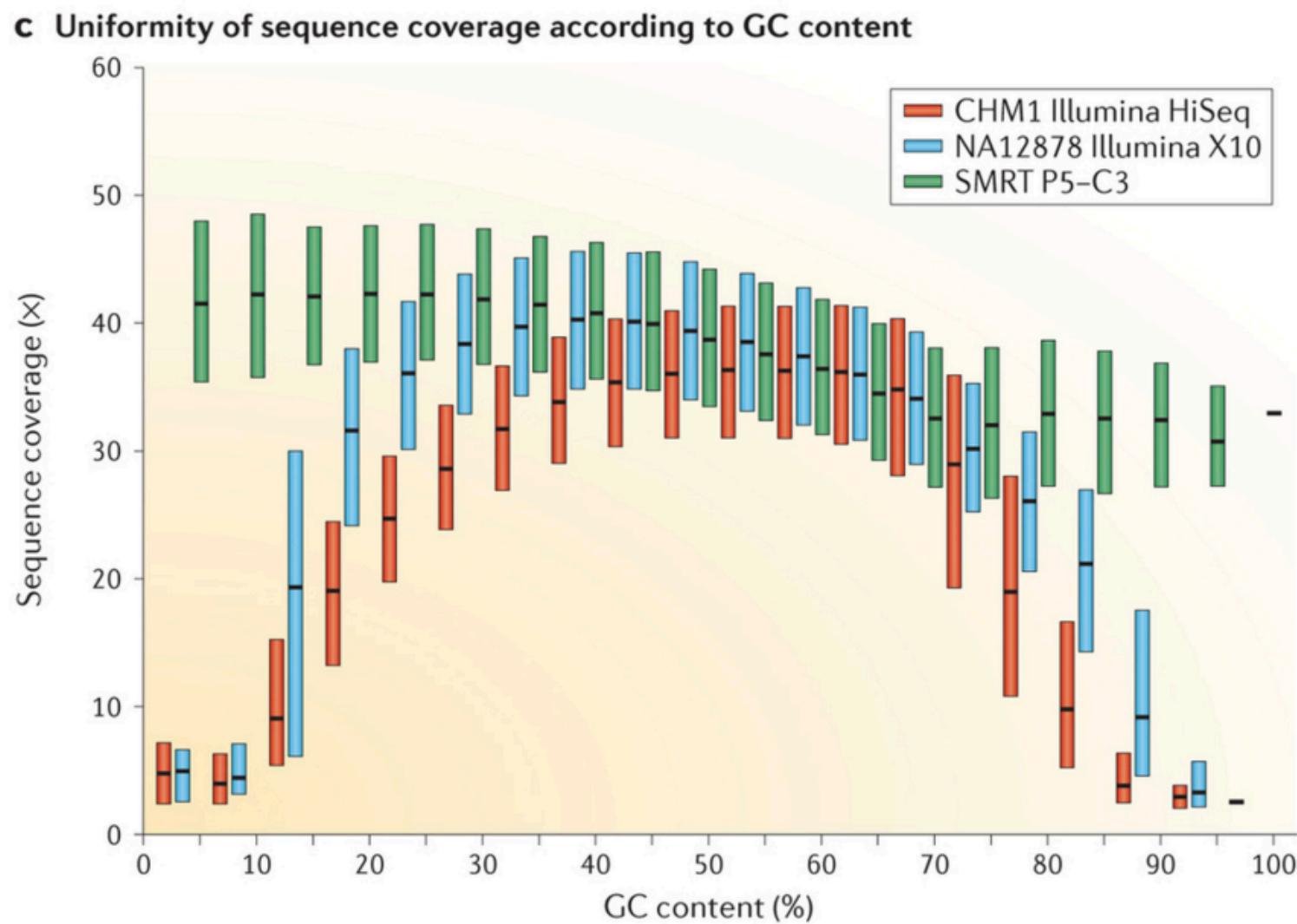


» Random entries



Analysis

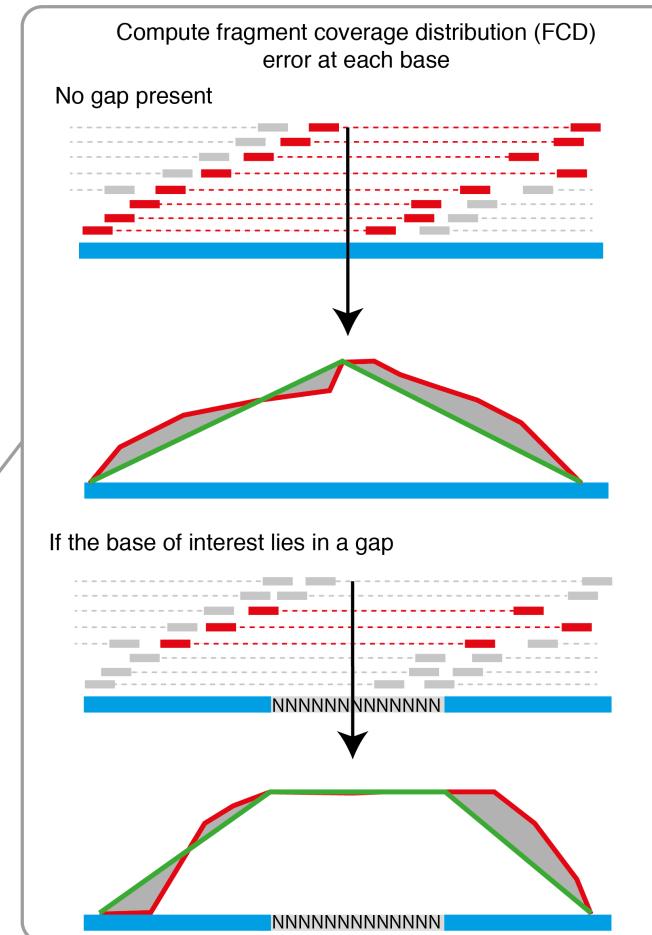
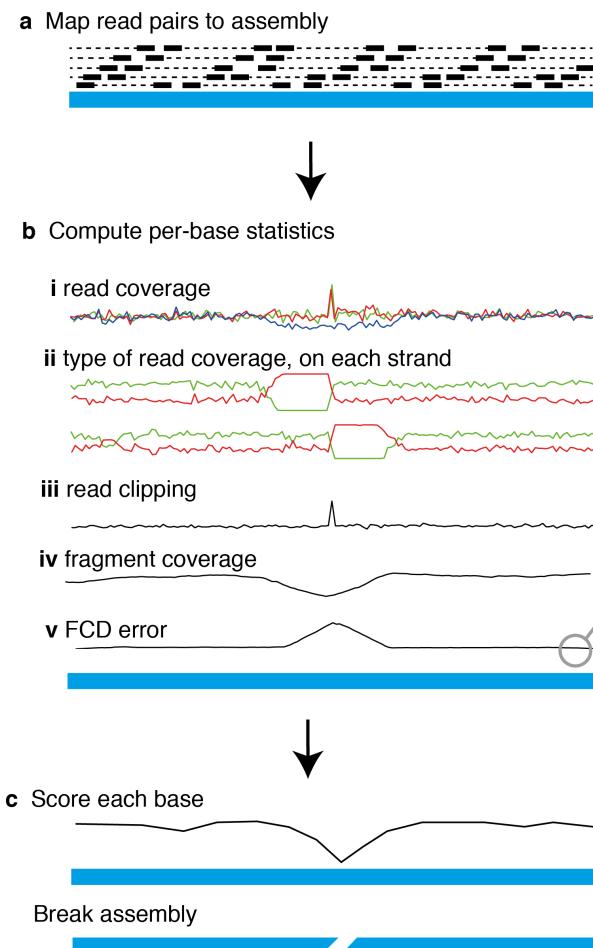
Biases



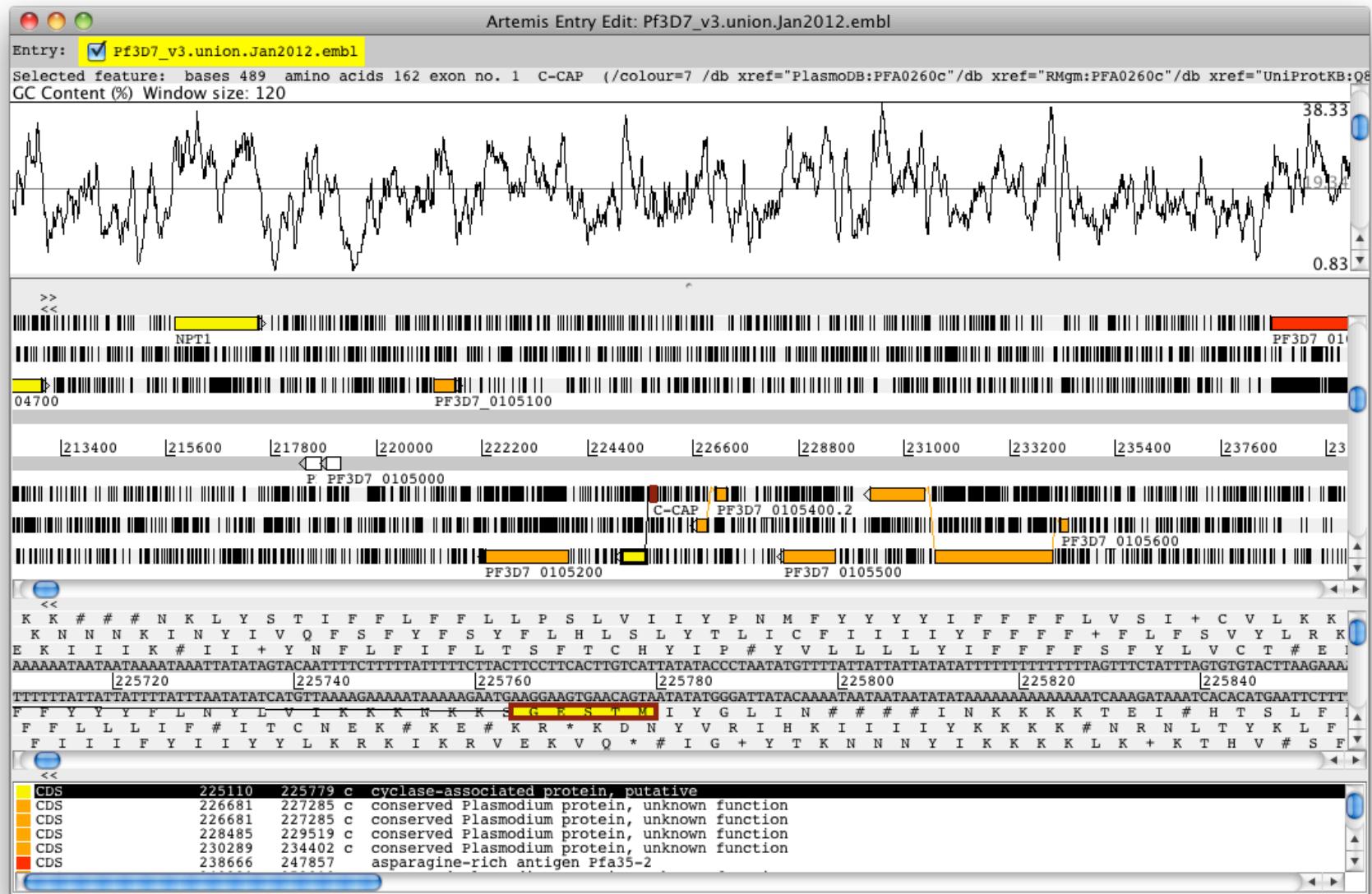
After assembly

- Say you have an assembly with 200 contigs and 34 scaffolds.
What do you do next?
- How accurate is it?
- Have you tried different assemblers?
- Can you improve with additional data or diminishing returns?
- Is there contamination?
- How does it compare to other species?

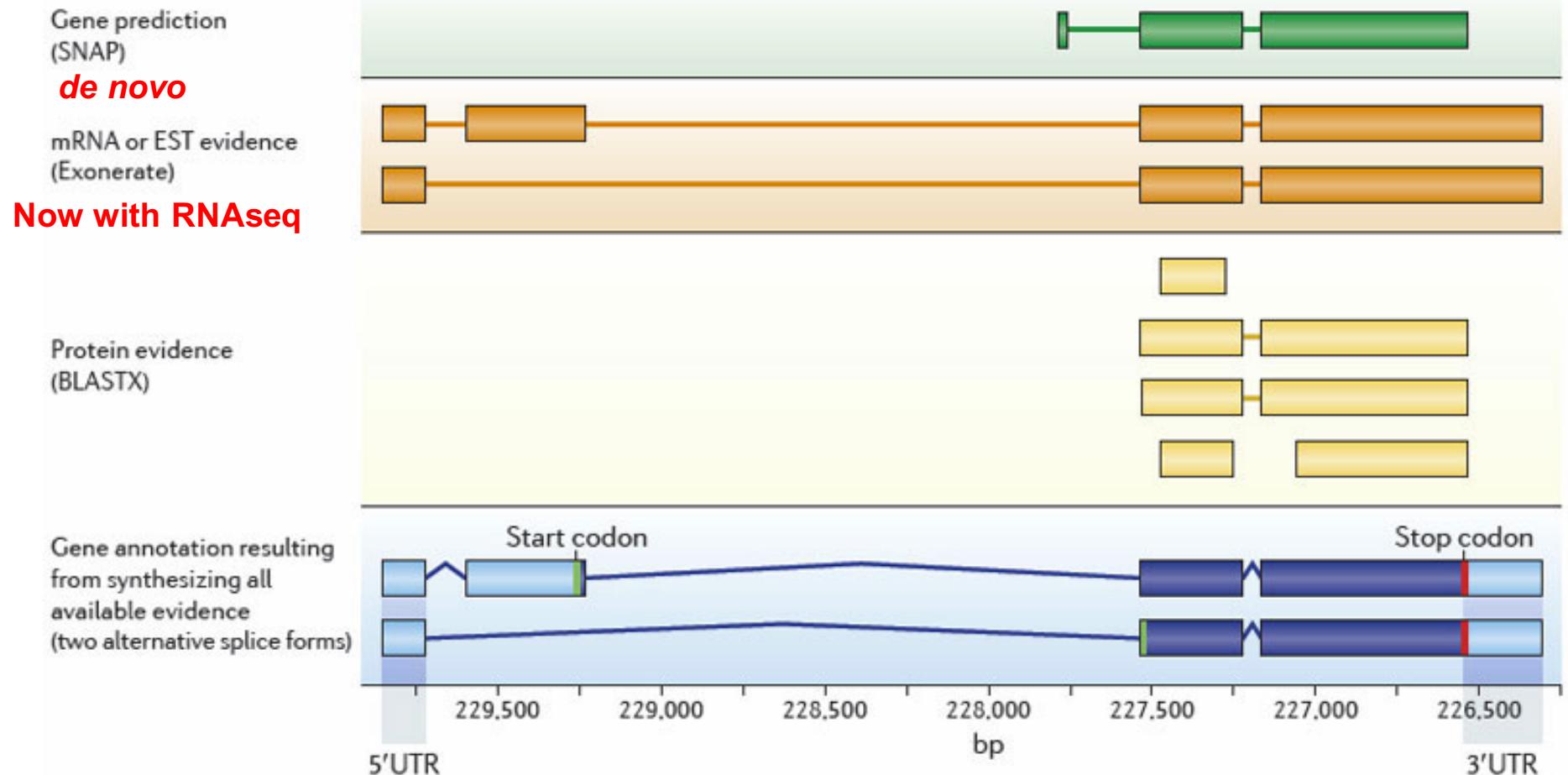
Reapr – Recognising Errors in Assemblies using Paired Reads



Visualisation (Artemis)



Annotation



Nature Reviews | Genetics

Yandell and Ence Nature Genetics Review (2012)

Classical genetics

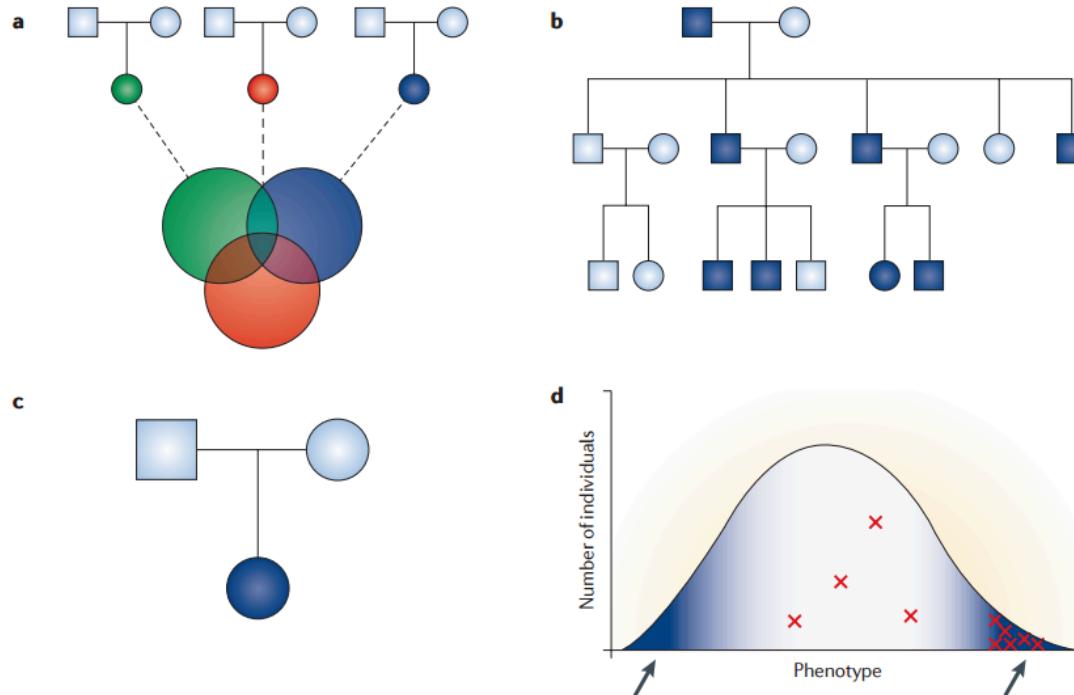
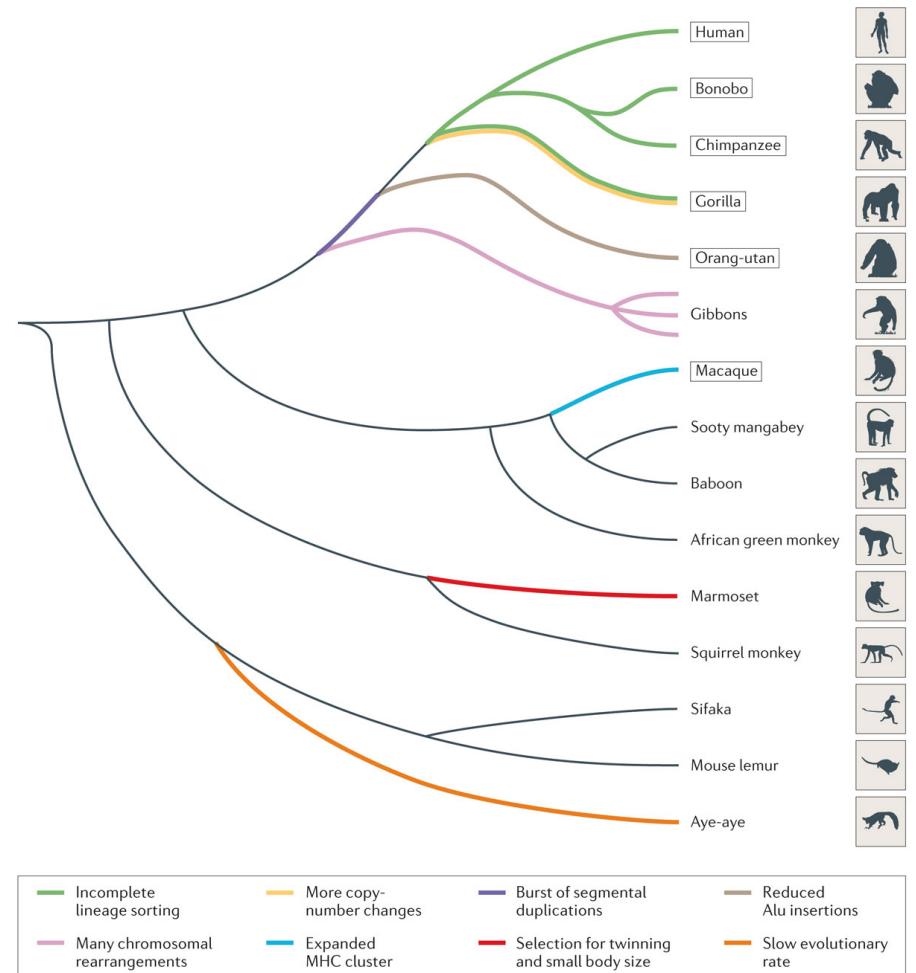
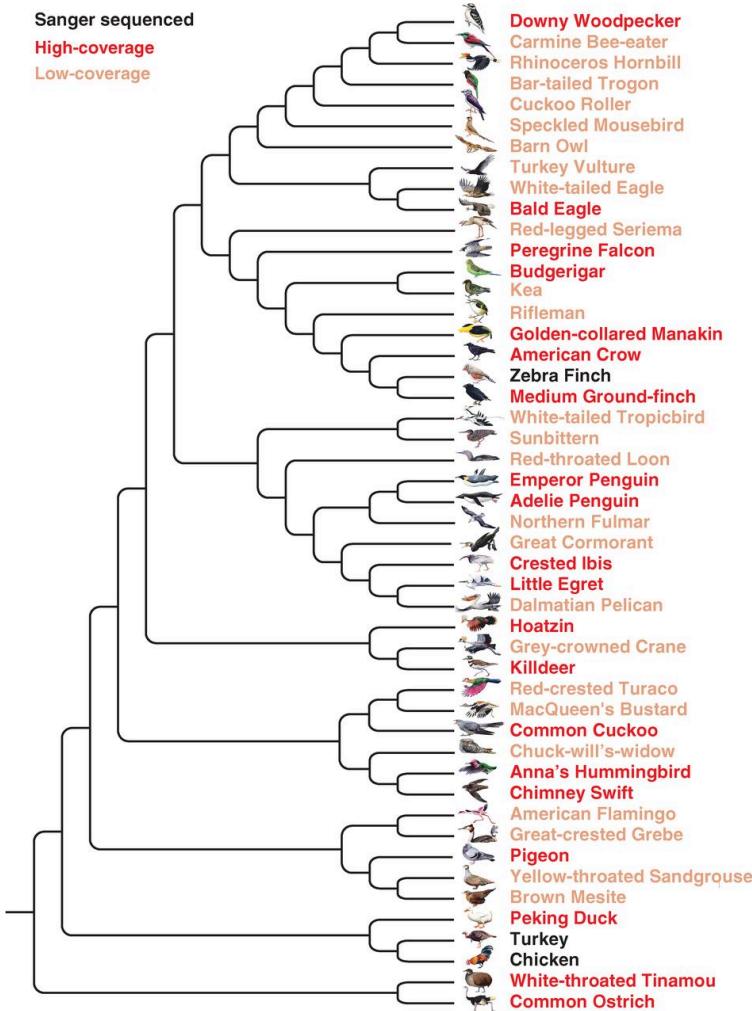


Figure 2 | Strategies for finding disease-causing rare variants using exome sequencing. Four main strategies are illustrated. **a** | Sequencing and filtering across multiple unrelated, affected individuals (indicated by the three coloured circles). This approach is used to identify novel variants in the same gene (or genes), as indicated by the shaded region that is shared by the three individuals in this example. **b** | Sequencing and filtering among multiple affected individuals from within a pedigree (shaded circles and squares) to identify a gene (or genes) with a novel variant in a shared region of the genome. **c** | Sequencing parent–child trios for identifying *de novo* mutations. **d** | Sampling and comparing the extremes of the distribution (arrows) for a quantitative phenotype. As shown in panel **d**, individuals with rare variants in the same gene (red crosses) are concentrated in one extreme of the distribution.

Comparative genomics



Nature Reviews | Genetics

Guojie Zhang et al. Science (2014)

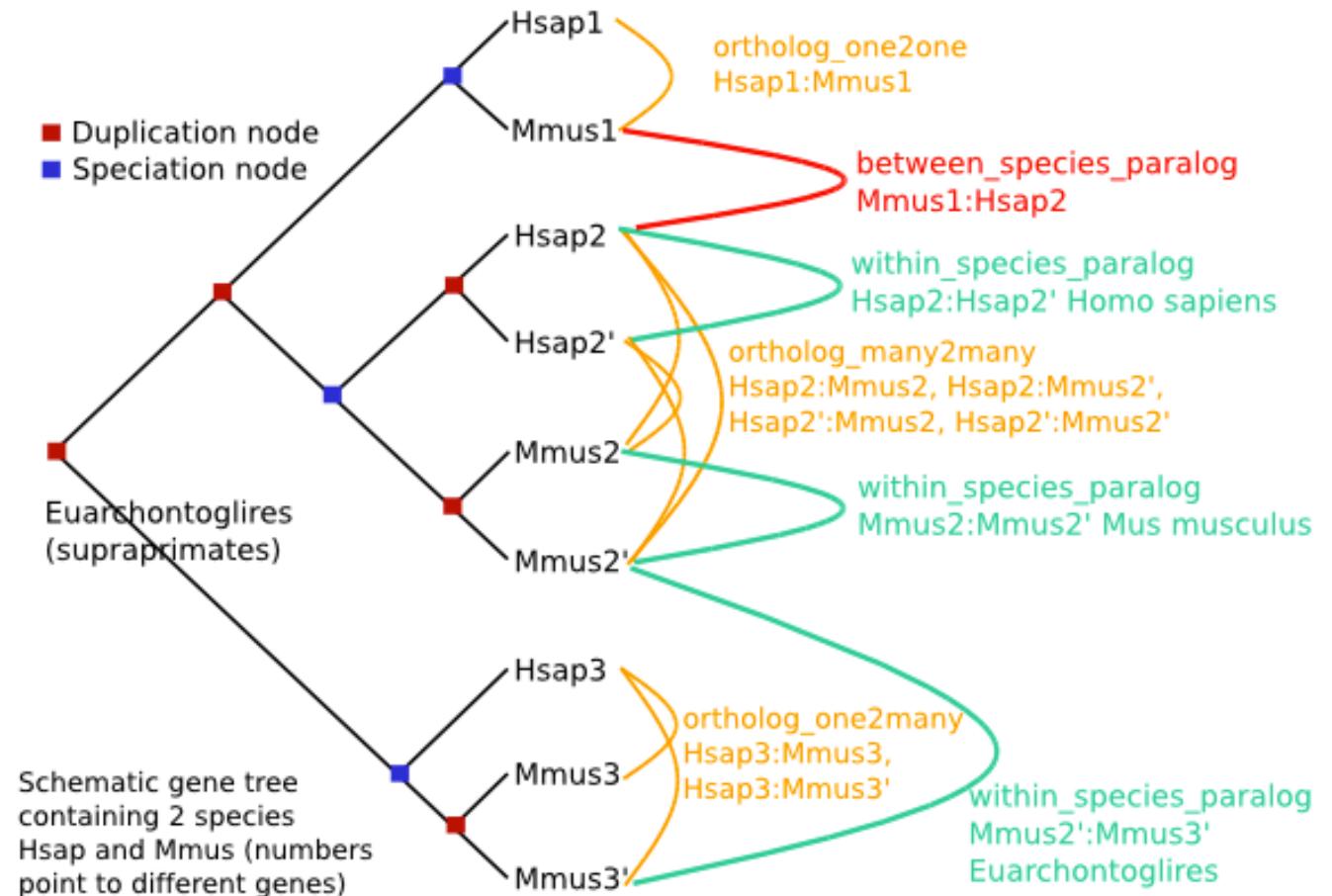
Roger & Gibbs Nature Reviews Genetics (2014)

Homologs: Orthologs and paralogs

Genes in different species and related by a speciation event are defined as **orthologs**.

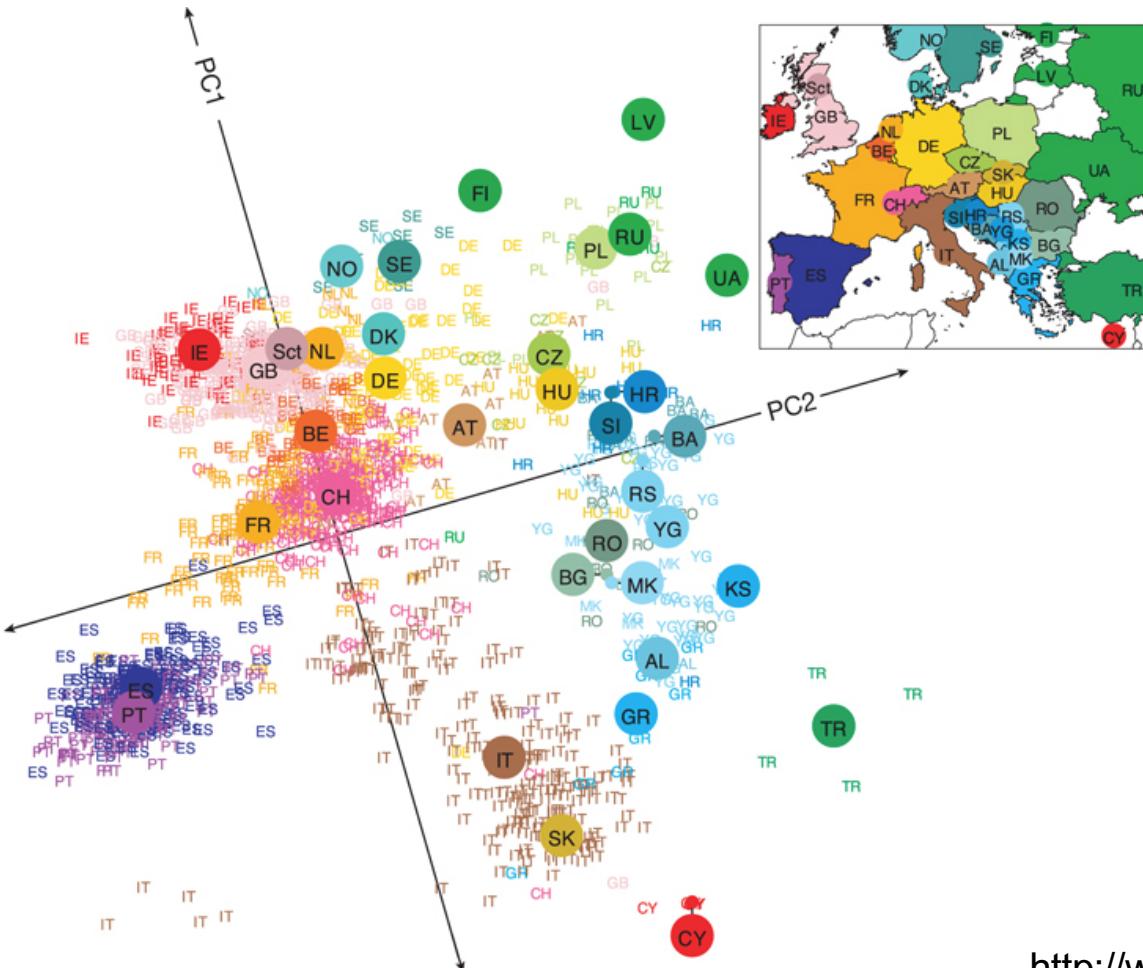
Depending on the number of genes found in each species, we differentiate among 1:1, 1:many and many:many relationships.

Genes of the same species and related by a duplication event are defined as **paralogs**.



http://asia.ensembl.org/info/genome/compara/homology_method.html?redirect=no

Population genomics

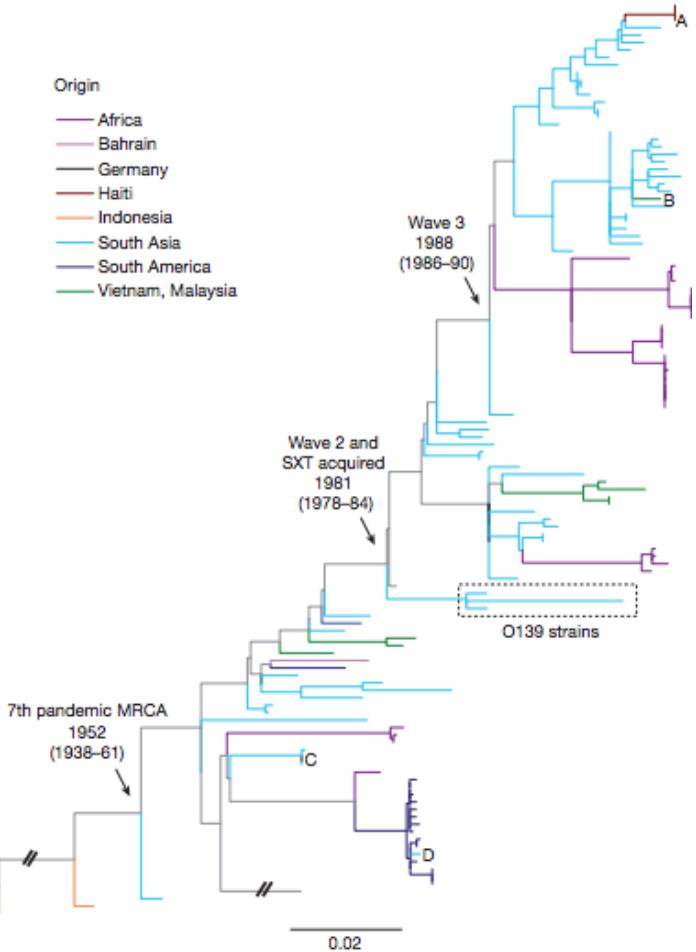
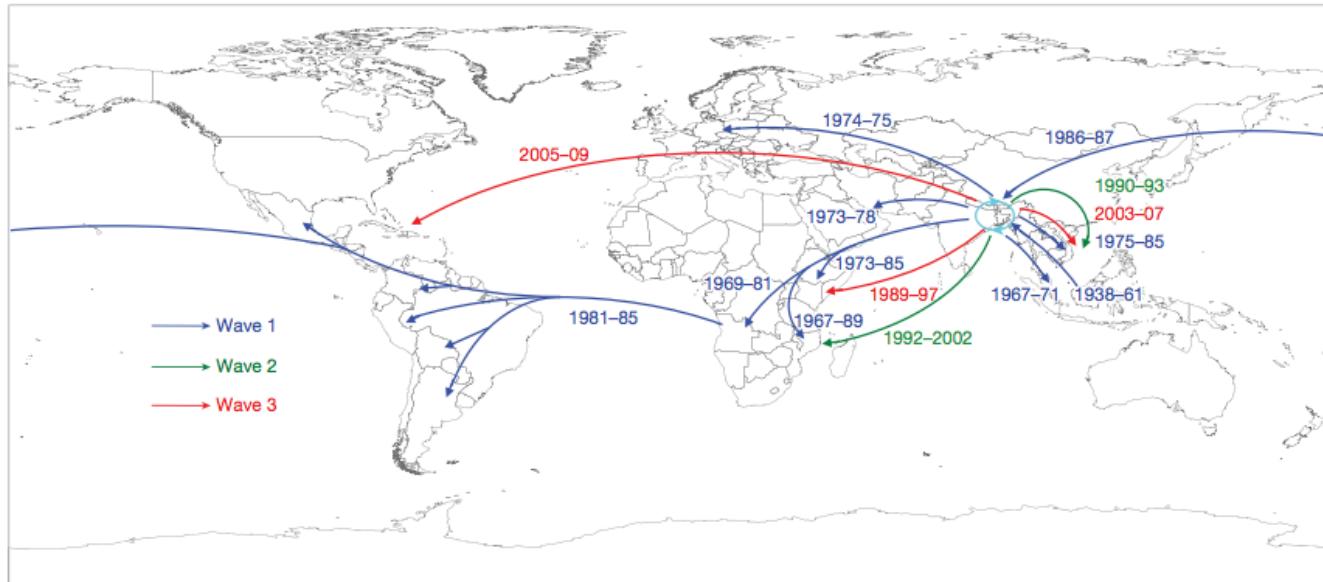


Novembre et al Nature (2008)



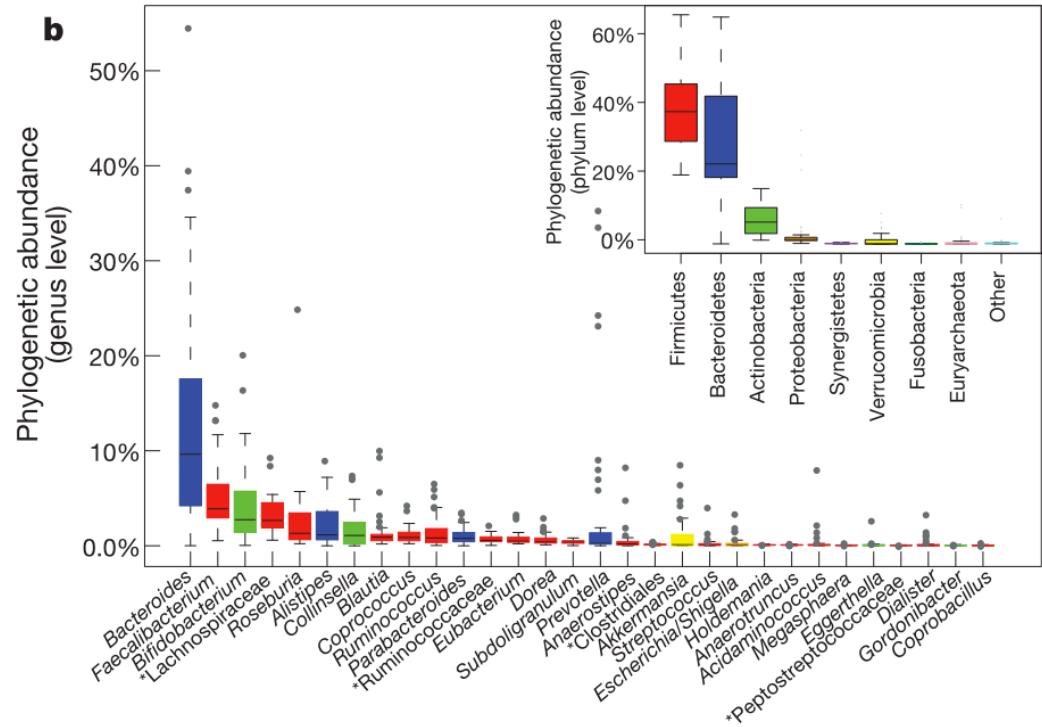
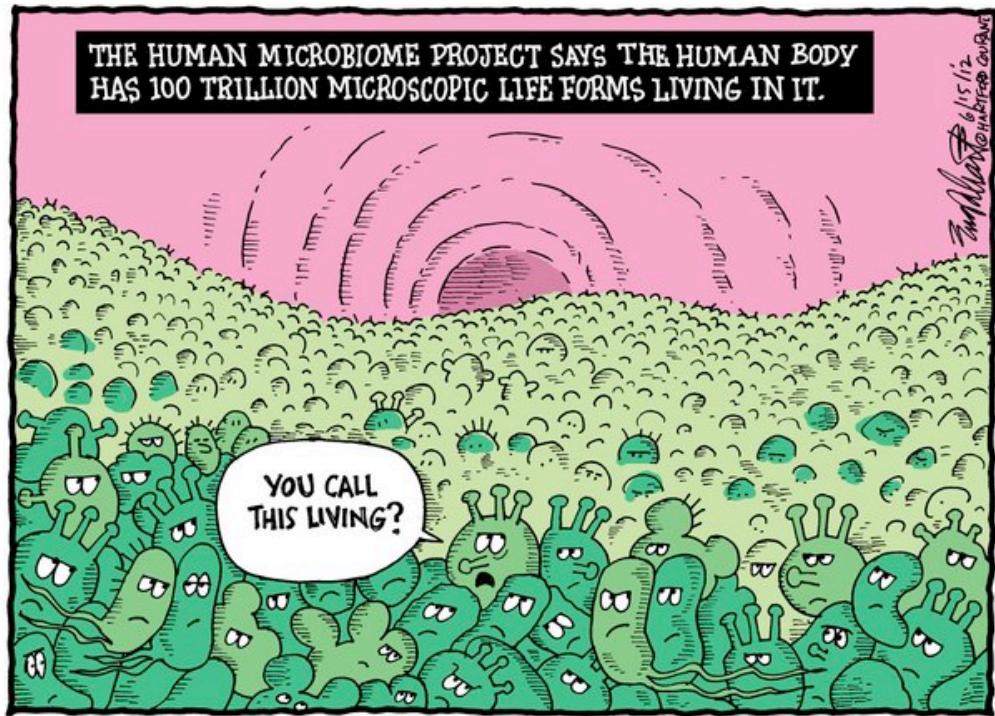
http://www.genomenext.com/casestudies_post/population-scale-analysis-genomic-samples-analyzed-from-2504-individuals-in-1-week/

Population genomics



Mutreja *et al.*, Nature Genetics (2011)

Metagenomics



doi:10.1038/nature09944

Personal journey

2005 – *Saccharomyces paradoxus*

- Capillary read sequenced full Chromosome III (~315kb) of 20 isolates
- Costed £750k
- One of the first scale re-sequencing projects
- Took me 3 years to sequence, align, annotate and analyse (= PhD)

Population genomics of the wild yeast *Saccharomyces paradoxus*: Quantifying the life cycle

Isheng J. Tsai, Douda Bensasson*, Austin Burt, and Vassiliki Koufopanou†

Division of Biology, Imperial College London, Silwood Park, Ascot, Berks SL5 7PY, United Kingdom

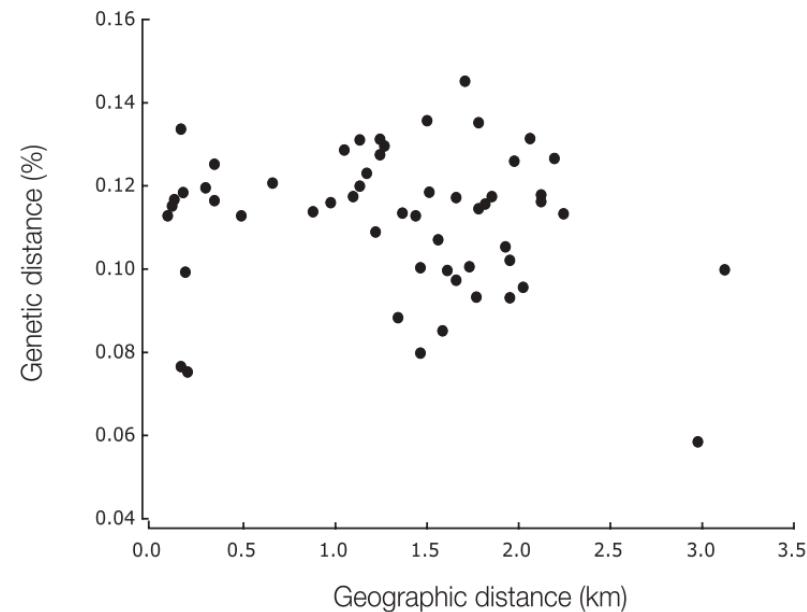
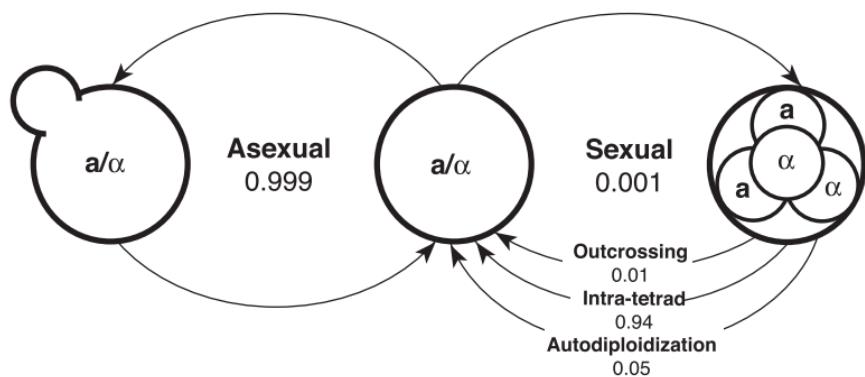
Edited by Mark Johnston, Washington University, St. Louis, MO, and accepted by the Editorial Board January 30, 2008 (received for review August 3, 2007)

Most microbes have complex life cycles with multiple modes of reproduction that differ in their effects on DNA sequence variation. Population genomic analyses can therefore be used to estimate the

are able to undergo mitoses, during which they repeatedly switch mating types, thus enabling matings between haploid clonemates (haplo-selfing or autodiploidization). This switch is possible be-

2005 – *Saccharomyces paradoxus*

- From population variation data we can infer frequencies of sex in yeast



2009 – *Saccharomyces* resequencing genome project

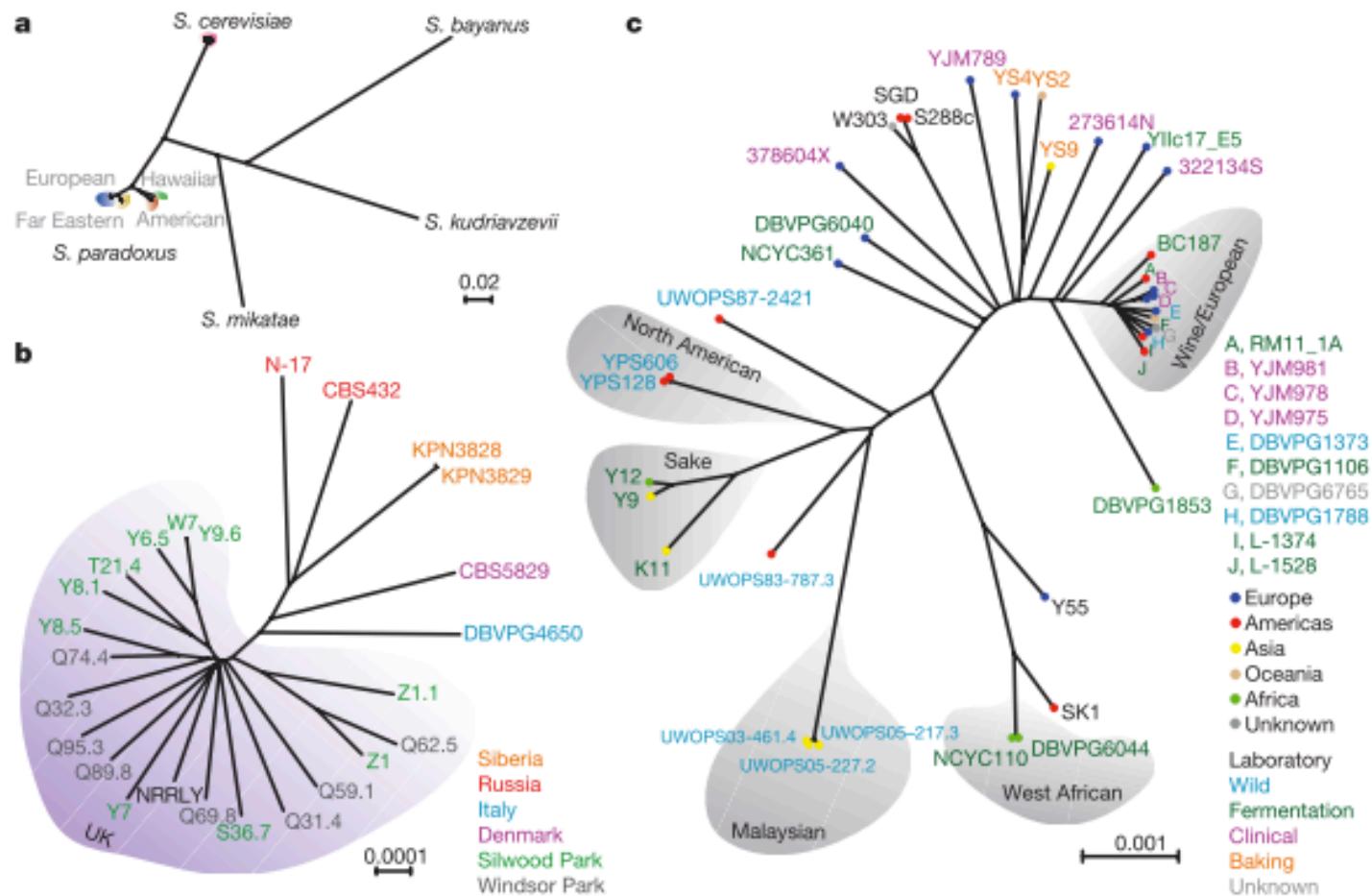
- 70 isolates at 1X-10X coverage
- ~2 years project with 26 authors
- At the start of NGS period (36bp Solexa reads)
- **Now = this can probably be done in 3 months by my RA**

Population genomics of domestic and wild yeasts

Gianni Liti^{1*}, David M. Carter^{2*}, Alan M. Moses^{2,3}, Jonas Warringer⁴, Leopold Parts², Stephen A. James⁵, Robert P. Davey⁵, Ian N. Roberts⁵, Austin Burt⁶, Vassiliki Koufopanou⁶, Isheng J. Tsai⁶, Casey M. Bergman⁷, Douda Bensasson⁷, Michael J. T. O'Kelly⁸, Alexander van Oudenaarden⁸, David B. H. Barton¹, Elizabeth Bailes¹, Alex N. Nguyen Ba³, Matthew Jones², Michael A. Quail², Ian Goodhead^{2†}, Sarah Sims², Frances Smith², Anders Blomberg⁴, Richard Durbin^{2*} & Edward J. Louis^{1*}

2009 – *Saccharomyces* resequencing genome project

Phylogeny of ~70 isolates



2013 – Tapeworm genome project

- 4 tapeworm genomes (~100Mb) of different sequencing technologies (Illumina, 454, capillary)
- RNAseq of host infecting cycle ; sequencing of 7 isolates
- 2 years of work with 56 authors

ARTICLE

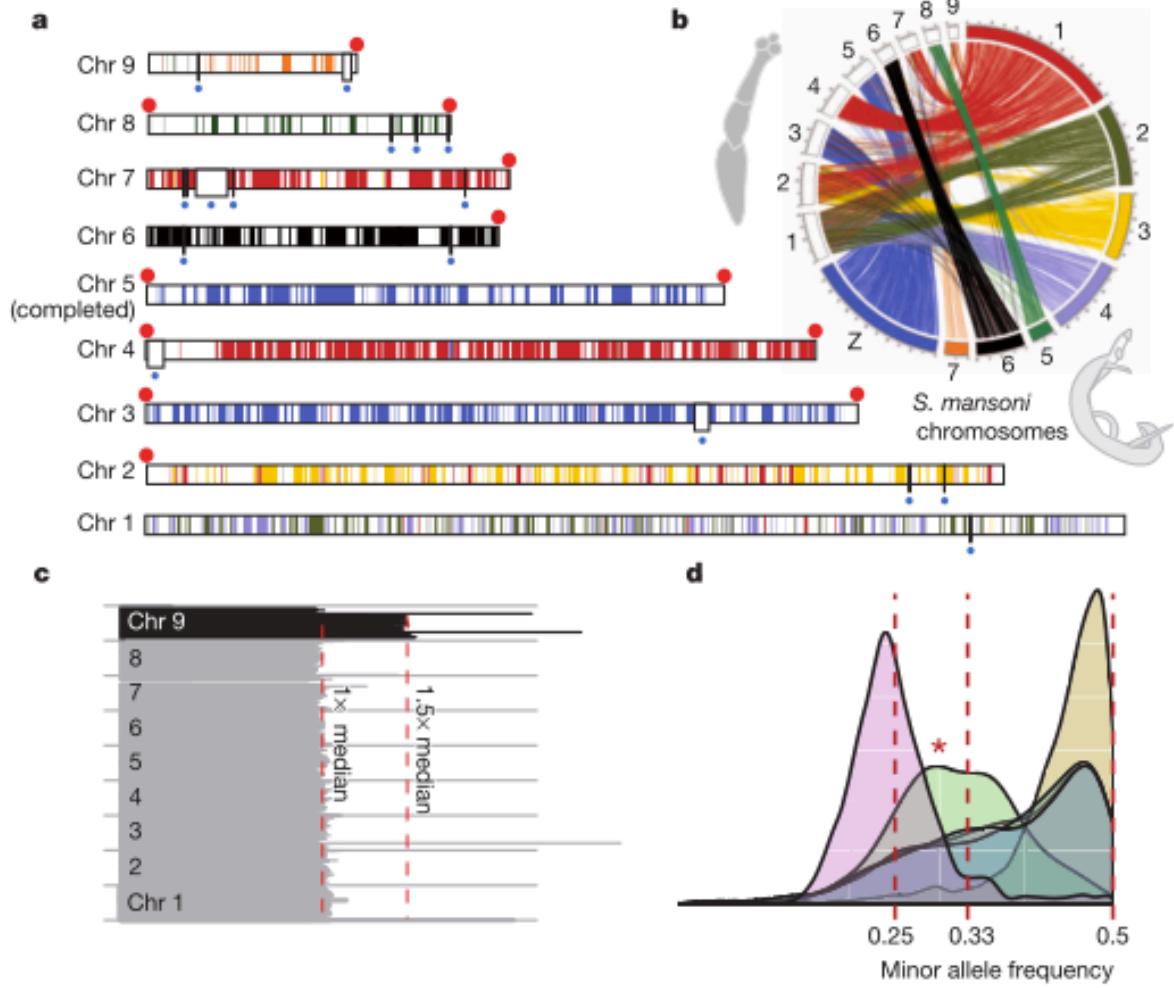
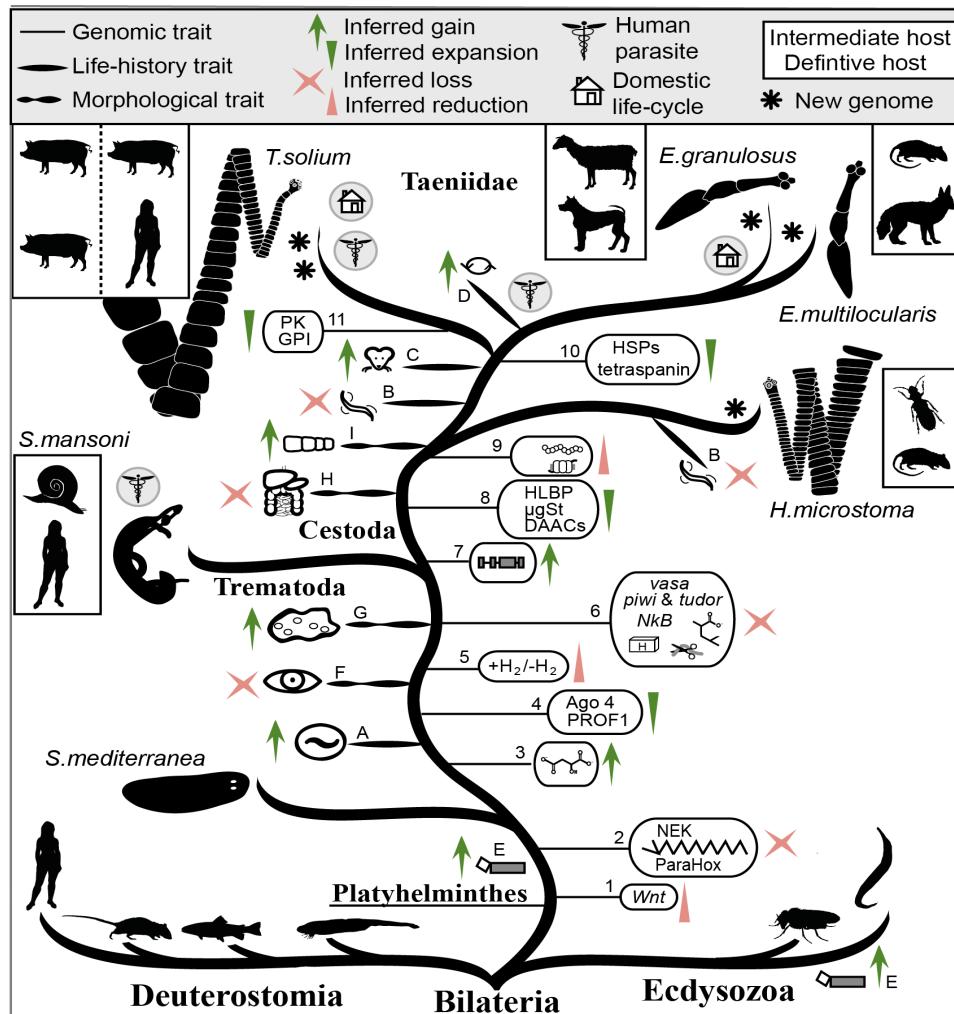
OPEN

doi:10.1038/nature12031

The genomes of four tapeworm species reveal adaptations to parasitism

Isheng J. Tsai^{1,2*}, Magdalena Zarowiecki^{1*}, Nancy Holroyd^{1*}, Alejandro Garciarrubio^{3*}, Alejandro Sanchez-Flores^{1,3}, Karen L. Brooks¹, Alan Tracey¹, Raúl J. Bobes⁴, Gladis Fragoso⁴, Edda Scuitto⁴, Martin Aslett¹, Helen Beasley¹, Hayley M. Bennett¹, Jianping Cai⁵, Federico Camicia⁶, Richard Clark¹, Marcela Cucher⁶, Nishadi De Silva¹, Tim A. Day⁷, Peter Deplazes⁸, Karel Estrada³, Cecilia Fernández⁹, Peter W. H. Holland¹⁰, Junling Hou⁵, Songnian Hu¹¹, Thomas Huckvale¹, Stacy S. Hung¹², Laura Kamenetzky⁶, Jacqueline A. Keane¹, Ferenc Kiss¹³, Uriel Koziol¹³, Olivia Lambert¹, Kan Liu¹¹, Xuenong Luo⁵, Yingfeng Luo¹¹, Natalia Macchiaroli⁶, Sarah Nichol¹, Jordi Paps¹⁰, John Parkinson¹², Natasha Pouchkina-Stantcheva¹⁴, Nick Riddiford^{14,15}, Mara Rosenzvit⁶, Gustavo Salinas⁹, James D. Wasmuth¹⁶, Mostafa Zamanian¹⁷, Yadong Zheng⁵, The *Taenia solium* Genome Consortium†, Xuepeng Cai⁵, Xavier Soberón^{3,18}, Peter D. Olson¹⁴, Juan P. Laclette⁴, Klaus Brehm¹³ & Matthew Berriman¹

2013 – Tapeworm genome project



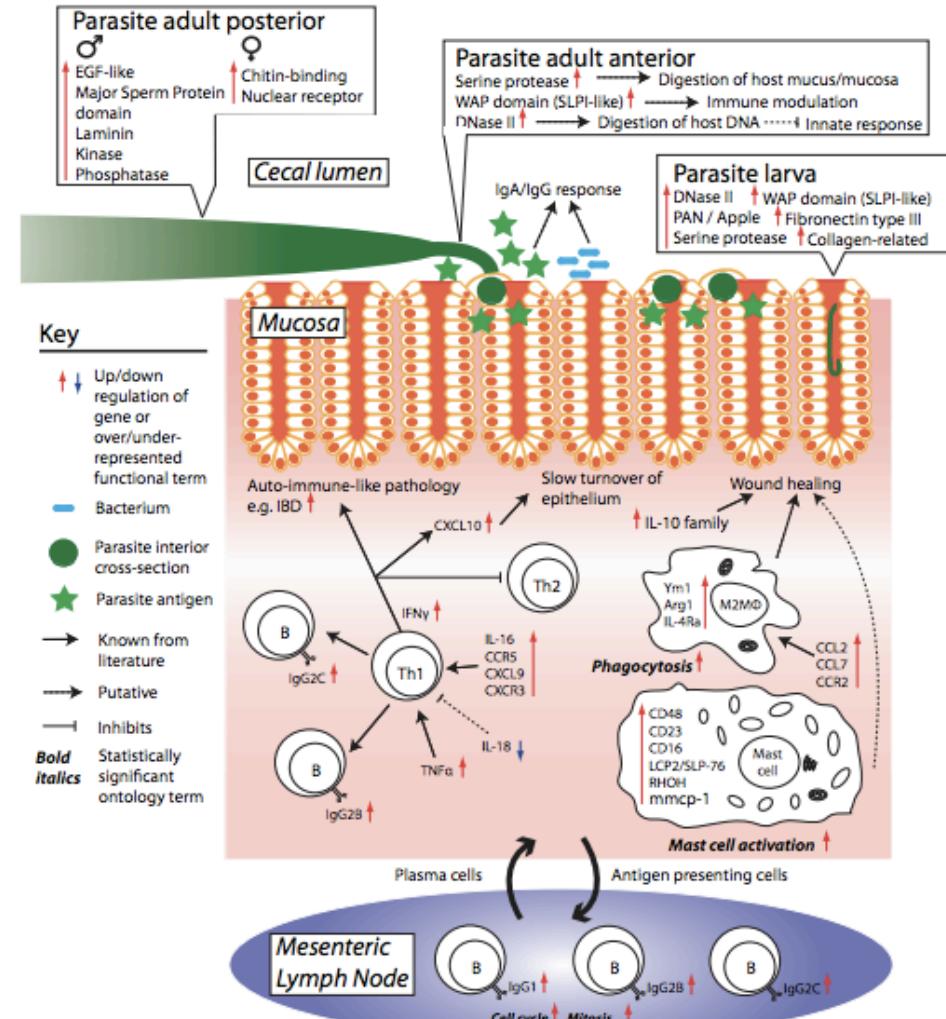
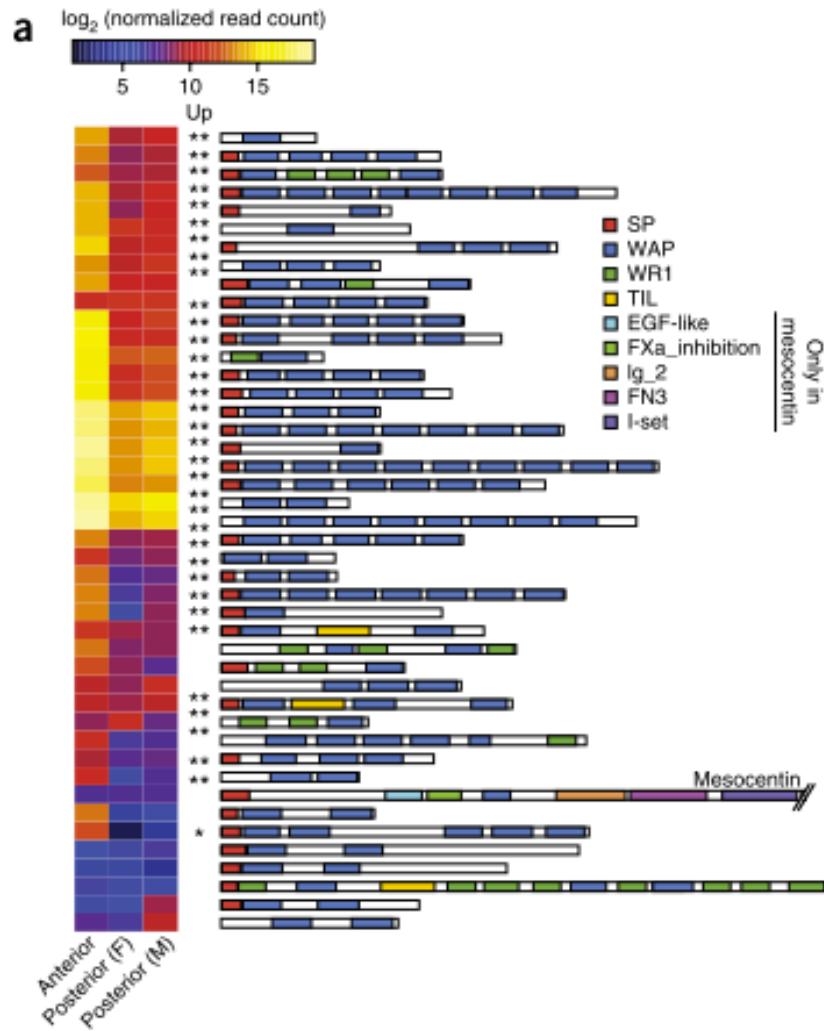
2014 – *Trichuris* genome project

- 2 genomes probably costs less than £10,000k
- About **40 RNAseq** libraries of different life cycle stages, host infecting stages
- Paradigm shifts to RNAseq

Whipworm genome and dual-species transcriptome analyses provide molecular insights into an intimate host-parasite interaction

Bernardo J Foth^{1,7}, Isheng J Tsai^{1,2,7}, Adam J Reid^{1,7}, Allison J Bancroft^{3,7}, Sarah Nichol¹, Alan Tracey¹, Nancy Holroyd¹, James A Cotton¹, Eleanor J Stanley¹, Magdalena Zarowiecki¹, Jimmy Z Liu⁴, Thomas Huckvale¹, Philip J Cooper^{5,6}, Richard K Grencis³ & Matthew Berriman¹

2014 – *Trichuris* genome project



2014 – *Taphrina* genome project

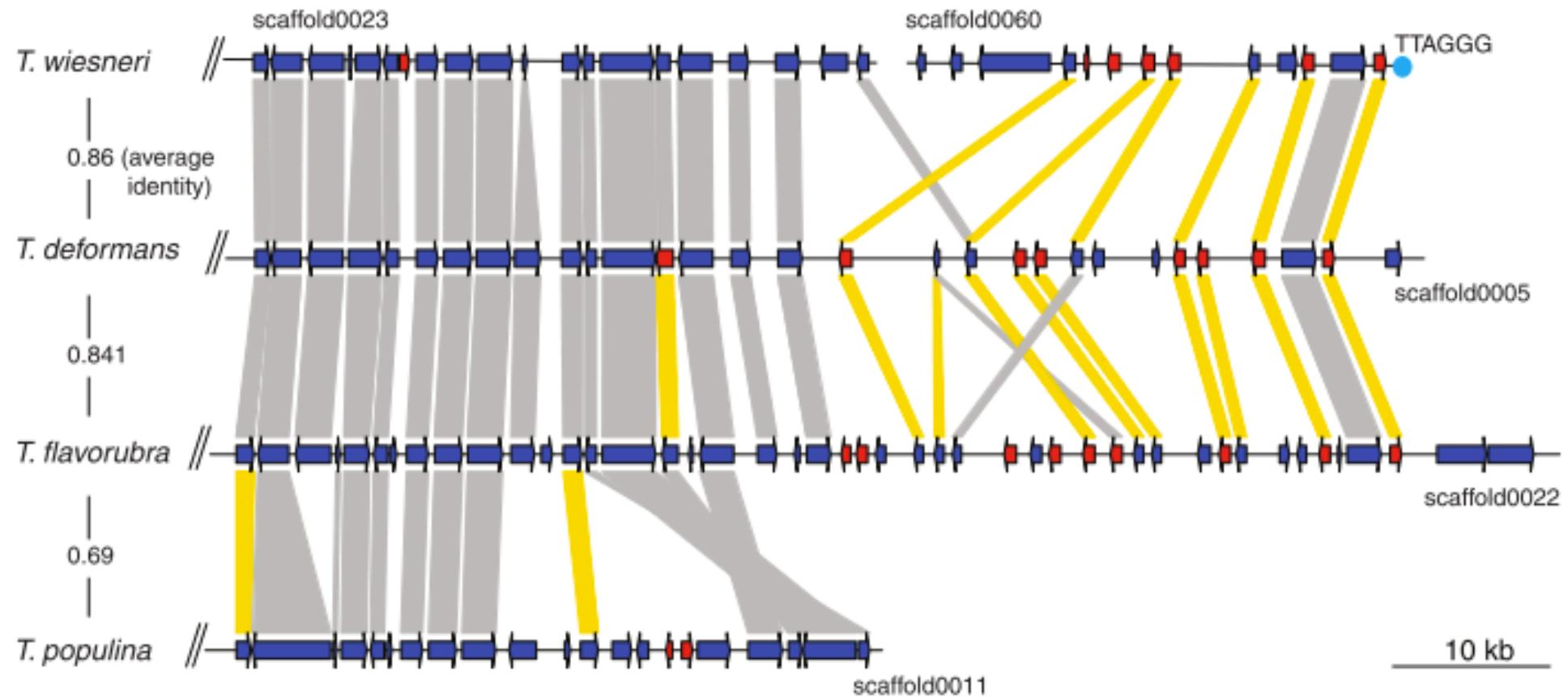
- 3 fungi genomes (~18Mb) of Illumina PE
- RNAseq for annotation purpose
- Costs probably less than 200,000 NT
- 2 months to analyse

GBE

Comparative Genomics of *Taphrina* Fungi Causing Varying Degrees of Tumorous Deformity in Plants

Isheng J. Tsai^{1,2}, Eiji Tanaka³, Hayato Masuya⁴, Ryusei Tanaka¹, Yuuri Hirooka⁵, Rikiya Endoh⁶, Norio Sahashi⁴, and Taisei Kikuchi^{1,4,*}

2014 – *Taphrina* genome project



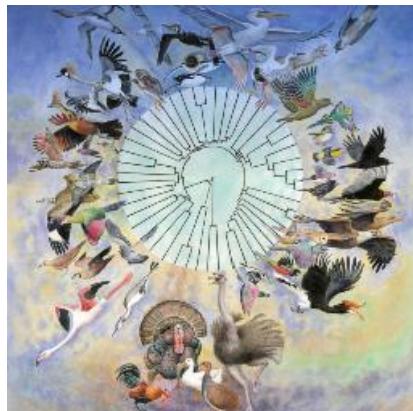
Shift in paradigm 2005-2015 (My personal take)

- A genome, a few genomes are no longer “enough”
 - ~since everybody can do it reasonably well
- Genome sequencing projects are being done on a per-lab basis
 - No longer exclusive to sequencing centres
 - But it also means some rubbish is being produced..
- Data being produced on a **much faster speed at a much higher throughput**, and a much **cheaper scale**
- More methods, analysis, tools, experiments...
 - Not always better

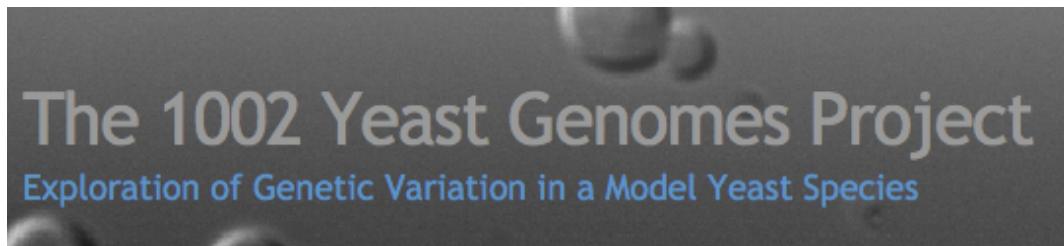
It is an exciting time to be in

Current and future

- Sequencing will still be cheaper, read will get longer
- Projects will be bigger



- Standard labs will be able to generate collections of themselves



(3 labs)

The end

- Please email to give feedbacks
 - Anything you want to know more in this lecture?
 - Anything I missed out?
- Any papers that you would like to recommend?
- Any topics that you would like focus?