

NGS lecture

Isheng Jason Tsai

Chang Gung University 2020
2019.03.26



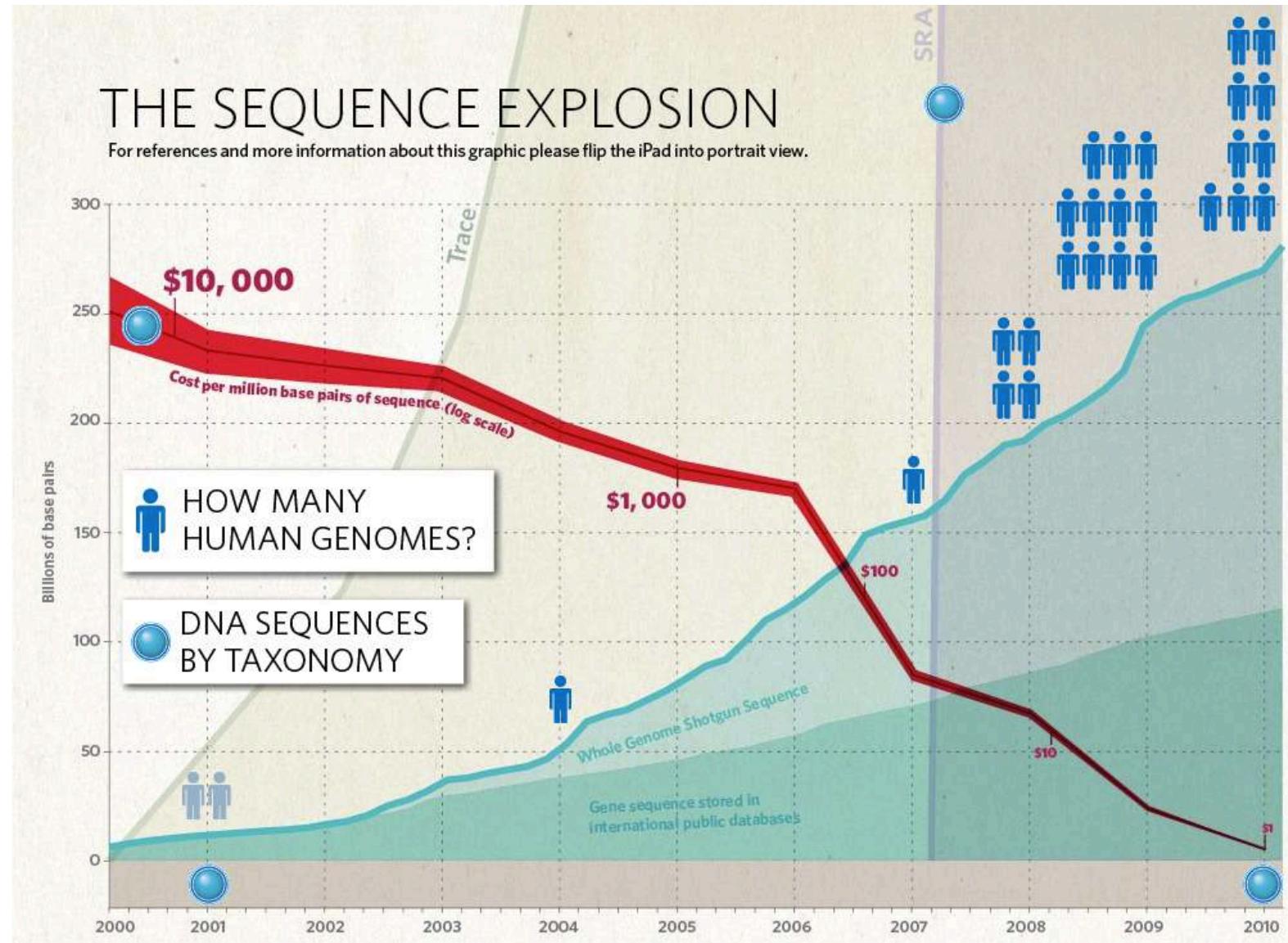
This lecture is called “NGS”

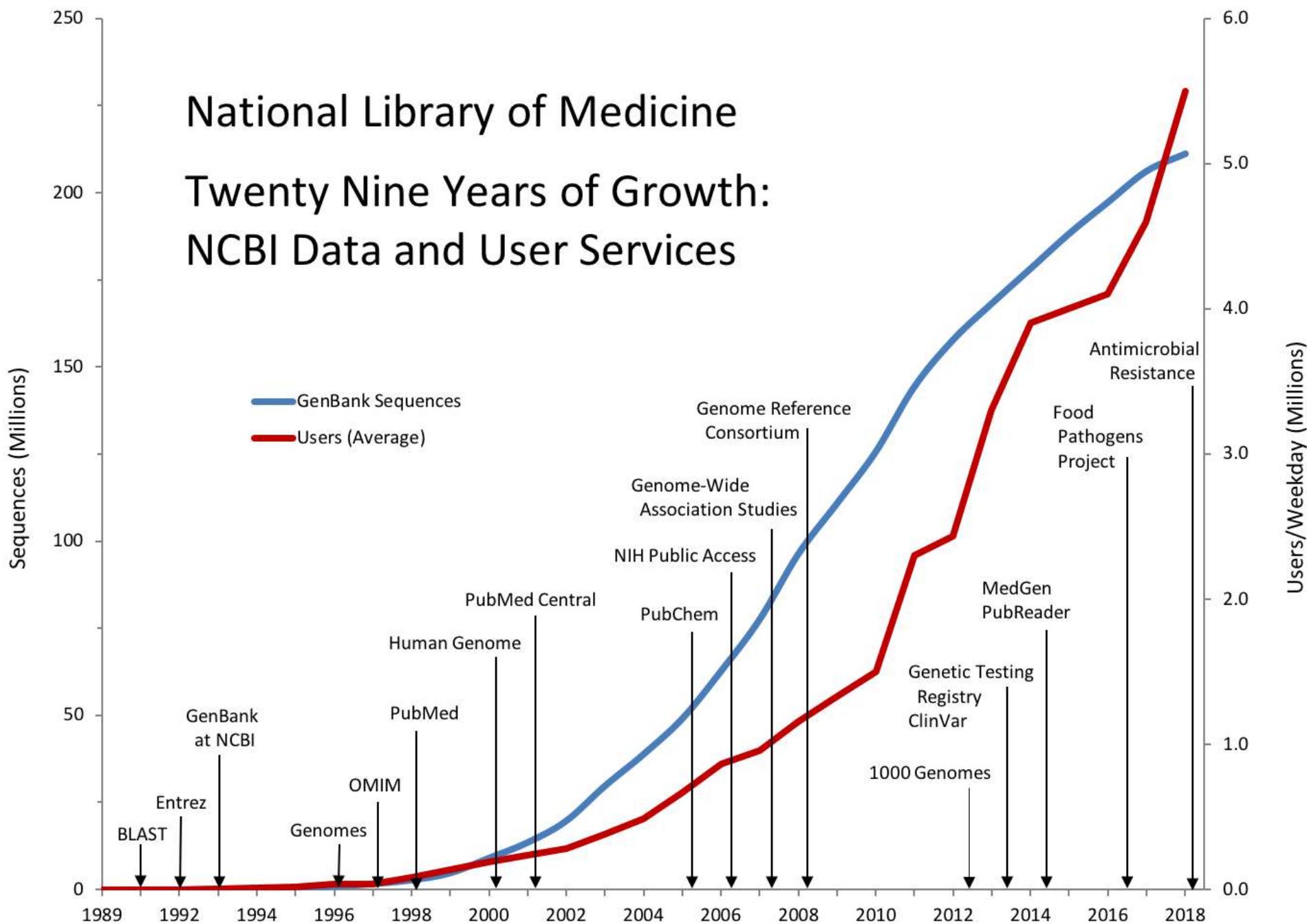
Actually

- 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?)

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





**NGS = sequencing made cheaper, faster and
higher throughput**

What we will cover today

NGS: Some basics

Sequencing platforms

Data types

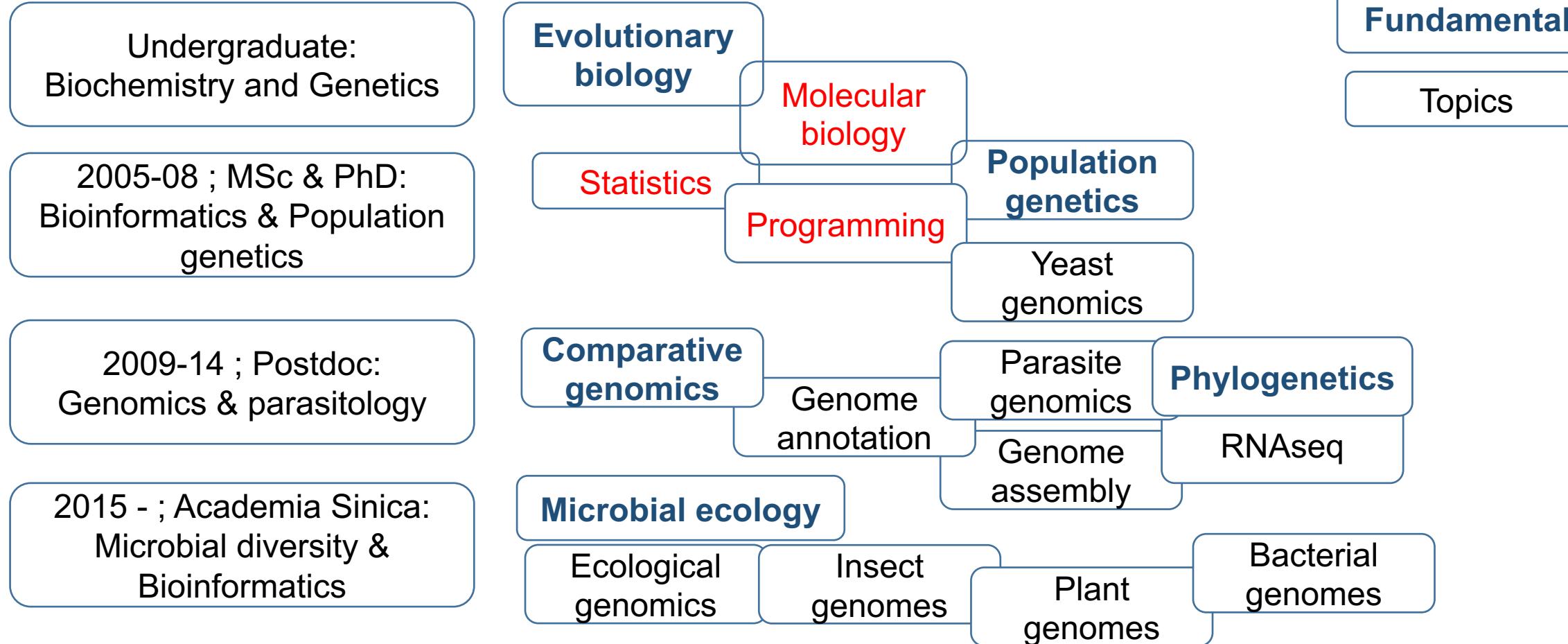
Analysis:

- RNAseq
- 16S
- Metagenomics

Previous questions:

- microbiota的paper要怎麼approach
- 16S sequencing region primer choice
- microbiota醫院有fecal transplantation計劃？

My background



- **49 publications**
(2 Nature, 1 Science, 2 Nature Genetics, 2 PNAS,
3 Genome Biology, 1 Nature Plant, 1 Nature Communications)

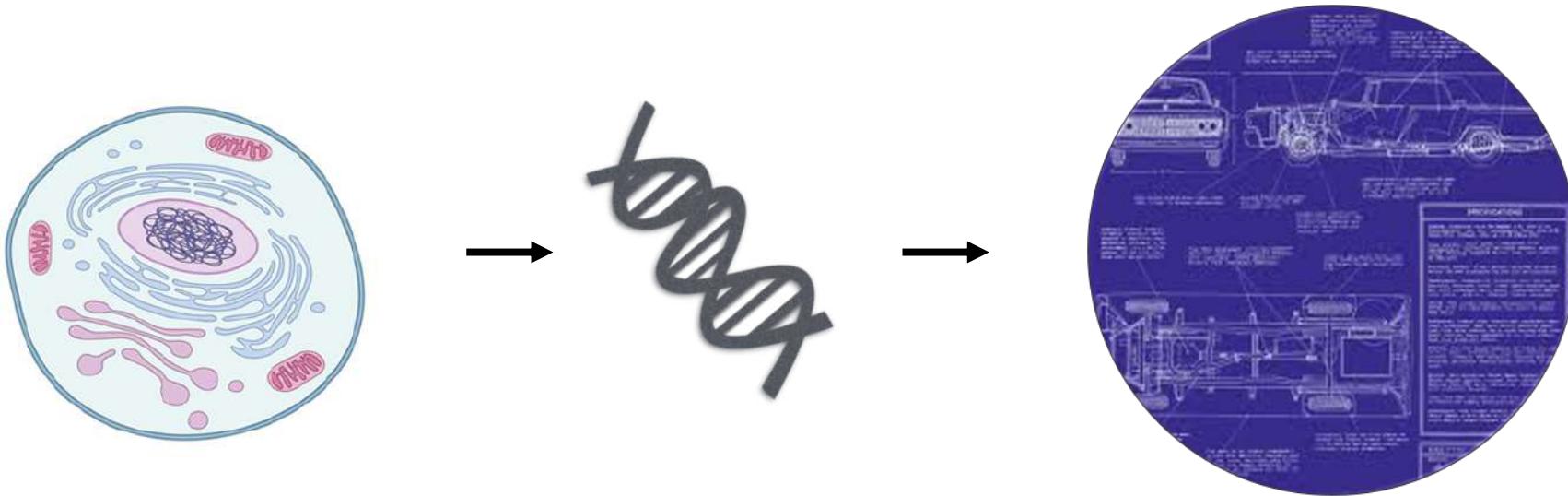
Skills

Fundamentals

Topics

What is genomics 基因體學？

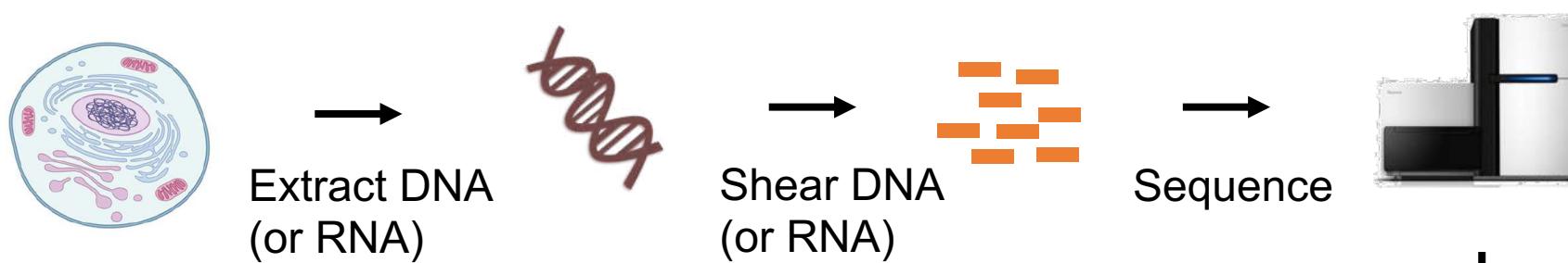
Genome



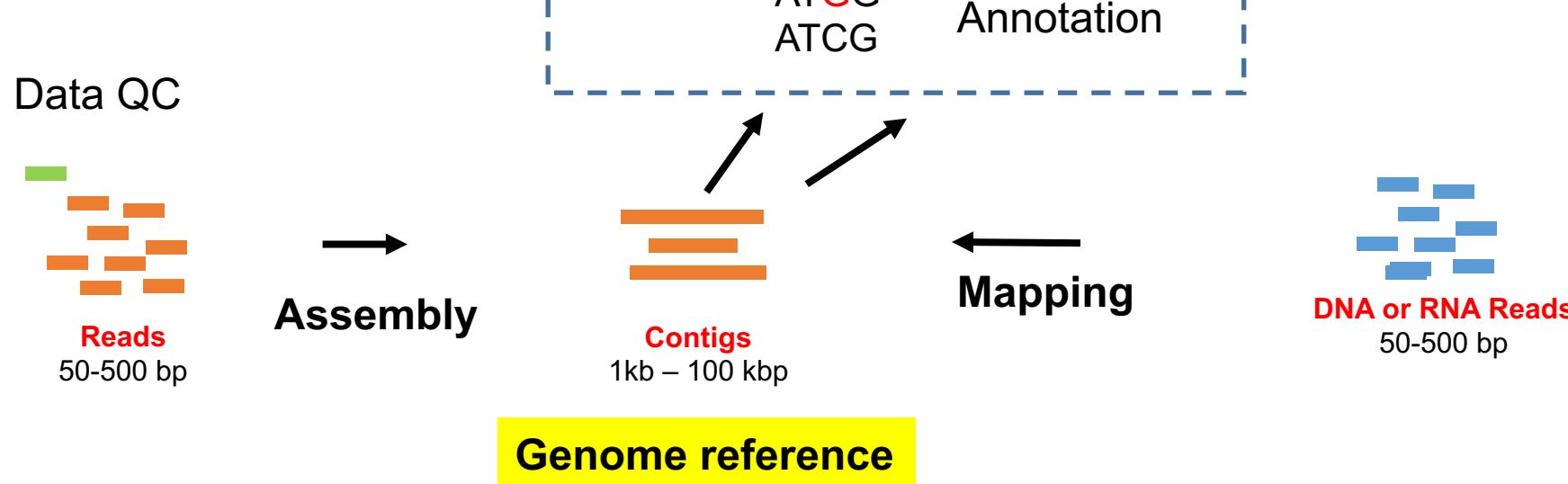
Genome = Parts list of a single genome

A genome project

Wet lab work



Bioinformatics



Three situations you are most likely to encounter

Genome reference is available (for example, humans):

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

Genome reference is NOT available

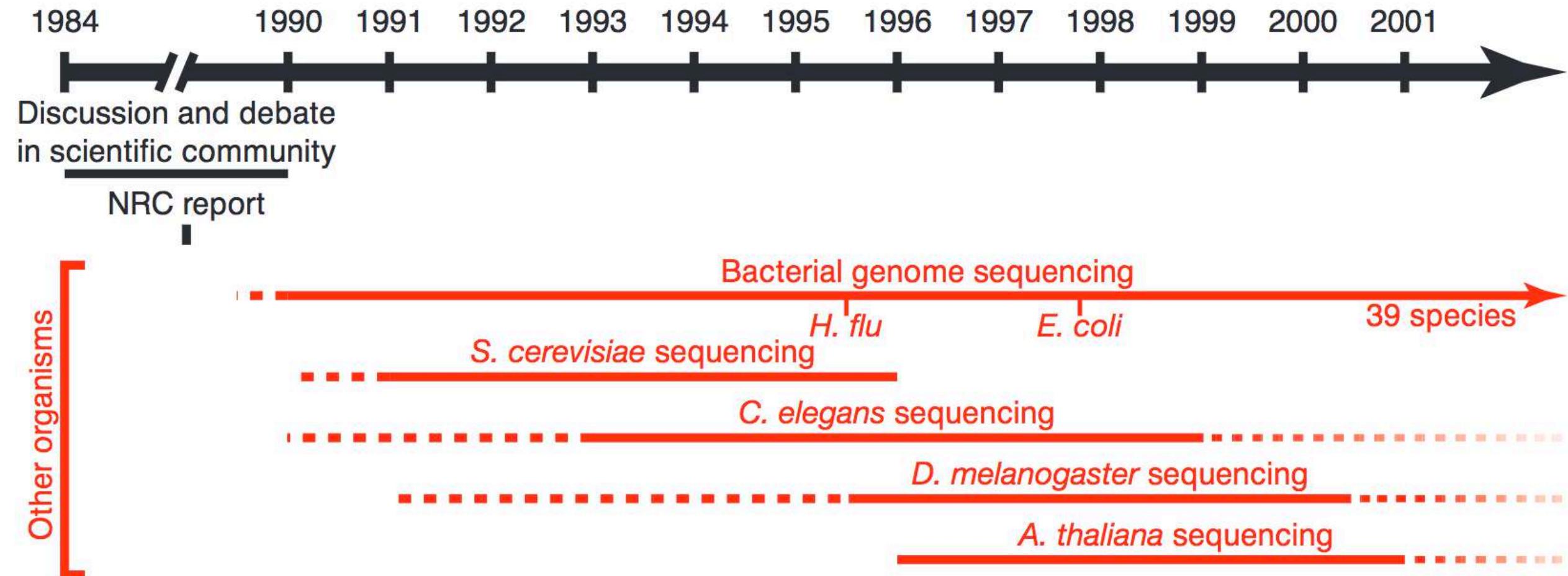
- **Assemble** the reads to get the genome

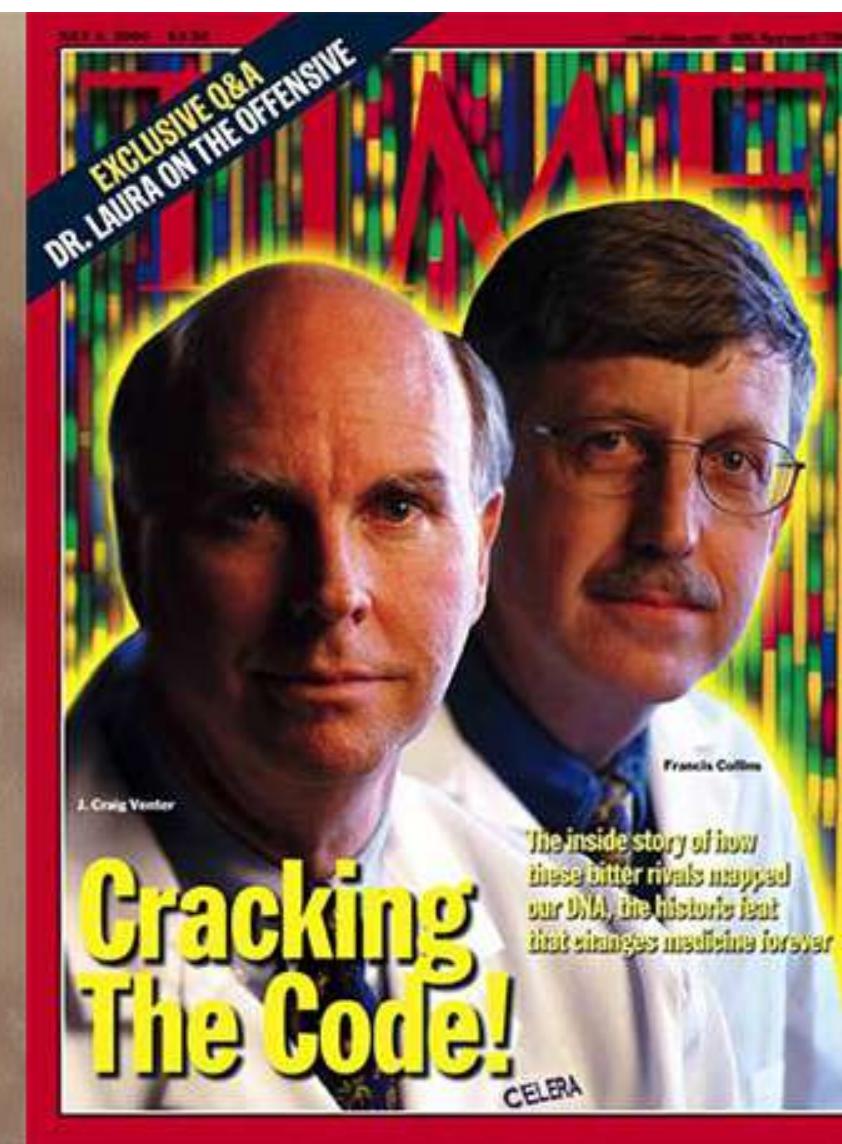
Counting:

- For a given region (gene) we want to know how much. → gene expression or metagenomics

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
- **Help to understand biology**





Calculating the economic impact of the Human Genome Project

Public funding of scientific R&D has a significant positive impact on the wider economy, but quantifying the exact impact of research can be difficult to assess. A new report by research firm Battelle Technology Partnership Practice estimates that **between 1988 and 2010, federal investment in genomic research generated an economic impact of \$796 billion**, which is impressive considering that Human Genome Project (HGP) spending **between 1990-2003 amounted to \$3.8 billion**. This figure equates to a return on investment (ROI) of 141:1 (that is, every \$1 invested by the U.S. government generated \$141 in economic activity). The report was commissioned by Life Technologies Foundation.

Large-scale whole-genome sequencing of the Icelandic population



A collection of Icelandic genealogical records dating back to the 1700s.

Here we describe the insights gained from sequencing the whole genomes of 2,636 Icelanders to a median depth of 20 \times .



The blood of a thousand Icelanders.
Photo: Chris Lund



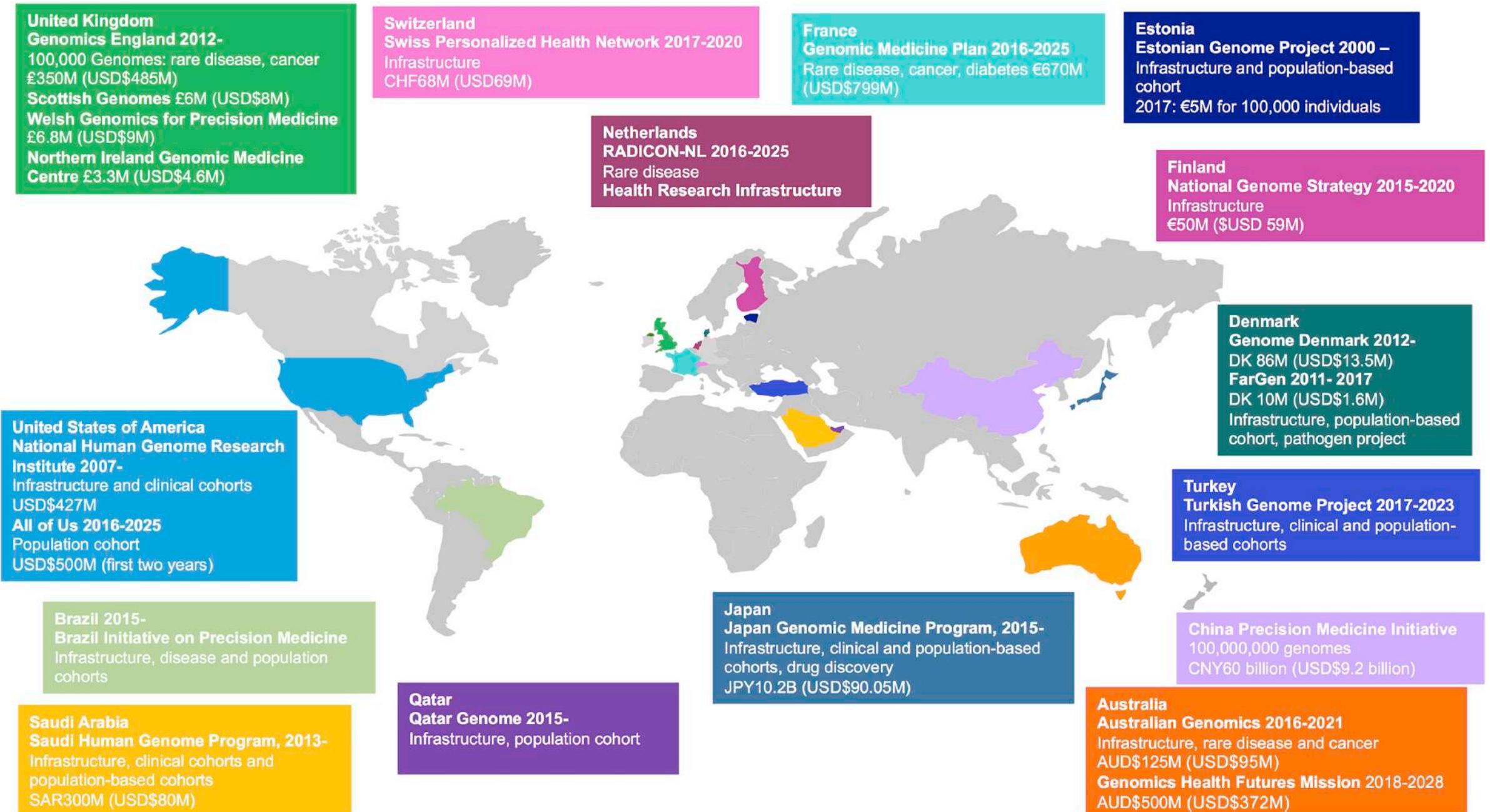
UK 10K

RARE GENETIC VARIANTS IN HEALTH AND DISEASE

The project is taking a two-pronged approach to identify rare variants and their effects:

- by studying and comparing the DNA of 4,000 people whose physical characteristics are well documented, the project aims to identify those changes that have no discernible effect and those that may be linked to a particular disease;
- by studying the changes within protein-coding areas of DNA that tell the body how to make proteins of 6,000 people with extreme health problems and comparing them with the first group, it is hoped to find only those changes in DNA that are responsible for the particular health problems observed.

The project received a £10.5 million funding award from Wellcome in March 2010 and sequencing started in late 2010. For more information, please use the links on the right hand side.





https://www.twbiobank.org.tw/new_web/index.php

The Cumulative 累計收案數

統計至2019年01月31日止([請按此](#))

社區民眾收案數

109,059

參與個案總數

22,502

完成第一輪追蹤個案總數

醫學中心患者收案數

1,862

參與個案總數

320

完成第一輪追蹤個案總數

8

完成第二輪追蹤個案總數

The Cumulative 累計收案數

統計至2020年02月29日止([請按此](#))

社區民眾收案數

127,853

參與個案總數

27,718

完成第一輪追蹤個案總數

醫學中心患者收案數

5,031

參與個案總數

882

完成第一輪追蹤個案總數

230

完成第二輪追蹤個案總數

21

完成第三輪追蹤個案總數

Whole exome sequencing and characterization of coding variation in 49,960 individuals in the UK Biobank

Authors: Cristopher V. Van Hout¹, Ioanna Tachmazidou², Joshua D. Backman¹, Joshua X. Hoffman², Bin Ye¹, Ashutosh K. Pandey², Claudia Gonzaga-Jauregui¹, Shareef Khalid¹, Daren Liu¹, Nilanjana Banerjee¹, Alexander H. Li¹, Colm O'Dushlaine¹, Anthony Marcketta¹, Jeffrey Staples¹, Claudia Schurmann¹, Alicia Hawes¹, Evan Maxwell¹, Leland Barnard¹, Alexander Lopez¹, John Penn¹, Lukas Habegger¹, Andrew L. Blumenfeld¹, Ashish Yadav¹, Kavita Praveen¹, Marcus Jones³, William J. Salerno¹, Wendy K. Chung⁴, Ida Surakka⁵, Cristen J. Willer⁵, Kristian Hveem⁶, Joseph B. Leader⁷, David J. Carey⁷, David H. Ledbetter⁷, Geisinger-Regeneron DiscovEHR Collaboration⁷, Lon Cardon², George D. Yancopoulos³, Aris Economides³, Giovanni Coppola¹, Alan R. Shuldiner¹, Suganthi Balasubramanian¹, Michael Cantor¹, Matthew R. Nelson^{2,*}, John Whittaker^{2,*}, Jeffrey G. Reid^{1,*}, Jonathan Marchini^{1,*}, John D. Overton^{1,*}, Robert A. Scott^{2,*}, Gonçalo Abecasis^{1,*}, Laura Yerges-Armstrong^{2,*}, Aris Baras^{1,*} on behalf of the Regeneron Genetics Center

The UK Biobank is a prospective study of 502,543 individuals, combining extensive phenotypic and genotypic data with streamlined access for researchers around the world.

	Variants in WES, n=49,960 Participants		Median Per Participant (IQR)	
	# Variants	# Variants MAF<1%	# Variants	# Variants MAF<1%
Total	9,693,526	9,547,730	48,982 (627)	1,626 (133)
Targeted Regions ¹	4,735,722	4,665,684	24,332 (283)	780 (63)
Variant Type¹				
SNVs	4,520,754	4,453,941	23,529 (276)	739 (61)
Indels	214,968	211,743	803 (29)	42 (10)
Multi-Allelic	591,340	580,728	3,388 (63)	117 (18)
Functional Prediction				
Synonymous	1,229,303	1,203,043	9,619 (128)	228 (28)
Missense	2,498,947	2,472,384	8,781 (137)	380 (39)
LOF (any transcript)	231,631	230,790	219 (16)	24 (8)
LOF (all transcripts)	153,903	153,441	111 (12)	15 (6)

Table 2 | Summary statistics for variants in sequenced exomes of 49,960 UKB participants

Project setup

- Sequencing a species (Comparative genomics)
 - Map, assemble
- **Sequencing multiple individuals of a species (Population genomics)**
 - **Map, count**
- Combination of (1) and (2)

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

生物資訊的開始

- the very beginnings of bioinformatics occurred **more than 50 years ago**, when desktop computers were still a hypothesis and DNA could not yet be sequenced.”
- The foundations of bioinformatics were laid in **the early 1960s** the application of computational methods to protein sequence analysis (notably, *de novo* sequence assembly, biological sequence databases and substitution models).
- Later on, DNA analysis also emerged due to parallel advances in (i) molecular biology methods, which allowed easier manipulation of DNA, as well as its sequencing, and (ii) computer science, which saw the rise of increasingly miniaturized and more powerful computers, as well as novel software better suited to handle bioinformatics tasks. **In the 1990s through the 2000s, major improvements in sequencing technology, along with reduced costs, gave rise to an exponential increase of data.**
- The arrival of ‘Big Data’ has laid out **new challenges in terms of data mining and management**, calling for more expertise from computer science into the field.

A brief history of bioinformatics

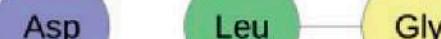
Jeff Gauthier, Antony T Vincent, Steve J Charette, Nicolas Derome
Briefings in Bioinformatics (2018) <https://doi.org/10.1093/bib/bby063>

**A****B**

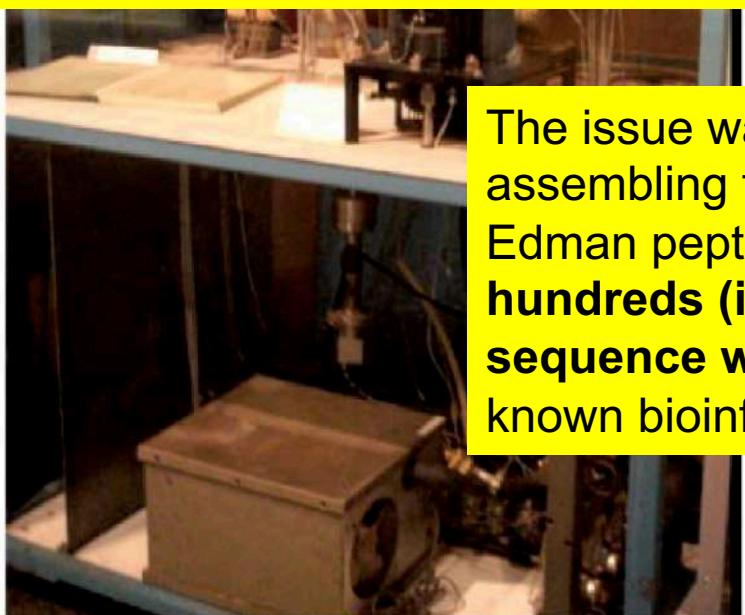
↓
Label 1st N-terminal residue



Remove 1st residue



A theoretical maximum of **50–60 amino acids** can be sequenced in a single Edman reaction. Larger proteins must be cleaved into smaller fragments, which are then separated and individually sequenced.



The issue was not sequencing a protein in itself but rather assembling the whole protein sequence from hundreds of small Edman peptide sequences. **For large proteins made of several hundreds (if not thousands) of residues, getting back the final sequence was cumbersome.** In the early 1960s, one of the first known bioinformatics software was developed to solve this problem.



Figure 1. Automated Edman peptide sequencing. (A) One of the first automated peptide sequencers, designed by William J. Dreyer. (B) Edman sequencing: the first N-terminal amino acid of a peptide chain is labeled with phenylisothiocyanate (PITC, red triangle), and then cleaved by lowering the pH. By repeating this process, one can determine a peptide sequence, one N-terminal amino acid at a time.

A brief history of bioinformatics

Jeff Gauthier, Antony T Vincent, Steve J Charette, Nicolas Derome
Briefings in Bioinformatics (2018) <https://doi.org/10.1093/bib/bby063>

Dayhoff: the first bioinformatician



Margaret Dayhoff (1925-1983)

- Designed one letter amino acid code
- Trained in quantum chemistry and mathematics, she became interested in proteins and molecular evolution around 1960.
- to explore mathematical approaches for analysing amino-acid sequence data
- Her initial project was writing a series of FORTRAN programs to determine the amino-acid sequences of protein molecules.

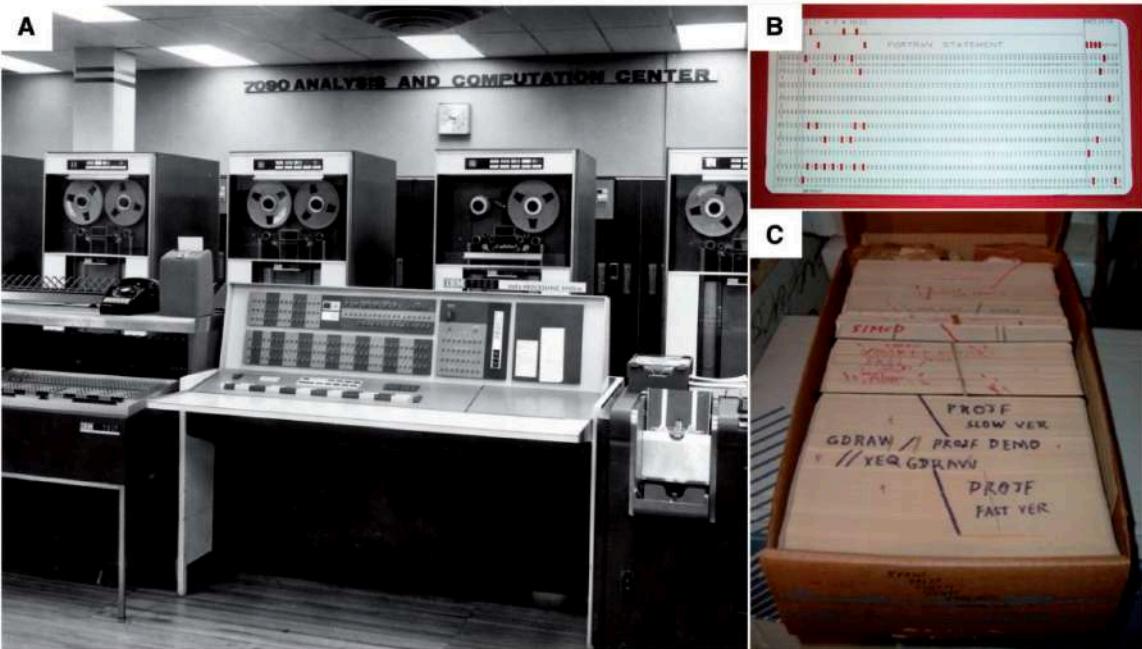
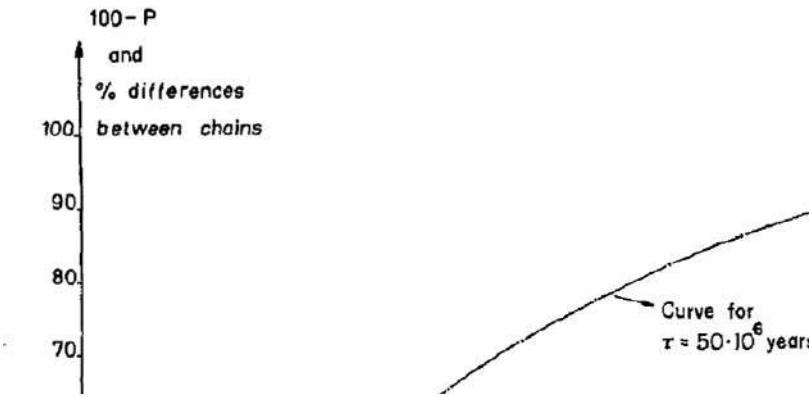
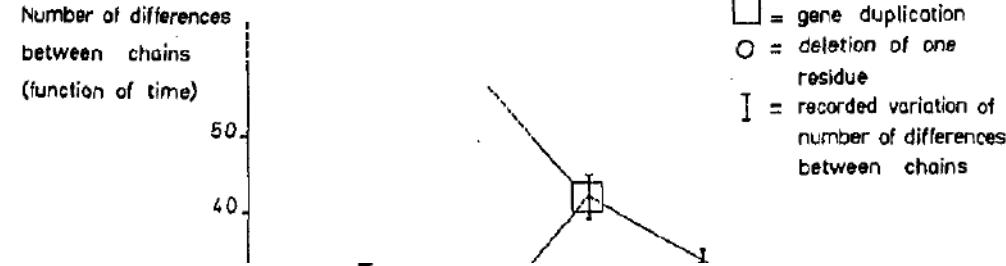


Figure 2. COMPROTEIN, the first bioinformatics software. (A) An IBM 7090 mainframe, for which COMPROTEIN was made to run. (B) A punch card containing one line of FORTRAN code (the language COMPROTEIN was written with). (C) An entire program's source code in punch cards. (D) A simplified overview of COMPROTEIN's input (i.e. Edman peptide sequences) and output (a consensus protein sequence).

A brief history of bioinformatics
Jeff Gauthier, Antony T Vincent, Steve J Charette, Nicolas Derome
Briefings in Bioinformatics (2018) <https://doi.org/10.1093/bib/bby063>

Ancestral sequences and Molecular clock (Emile Zuckerkandl and Linus Pauling)



There may thus exist a molecular evolutionary clock.

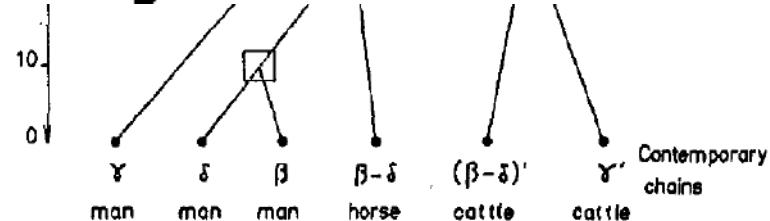
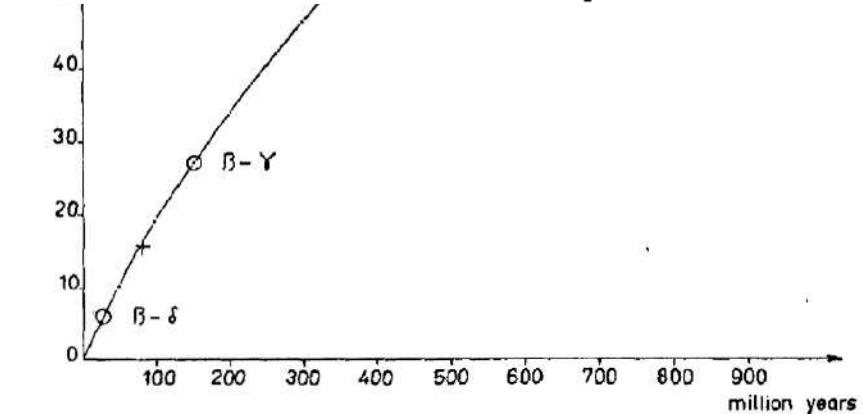


FIG. 4. Probable evolutionary relationship of some mammalian hemoglobin chains.

"Zuckerkandl and Pauling hypothesized that orthologous proteins evolved through divergence from a common ancestor. Consequently, by comparing the sequence of hemoglobin in currently extant organisms, it became possible to predict the 'ancestral sequences' of hemoglobin and, in the process, its evolutionary history up to its current forms"

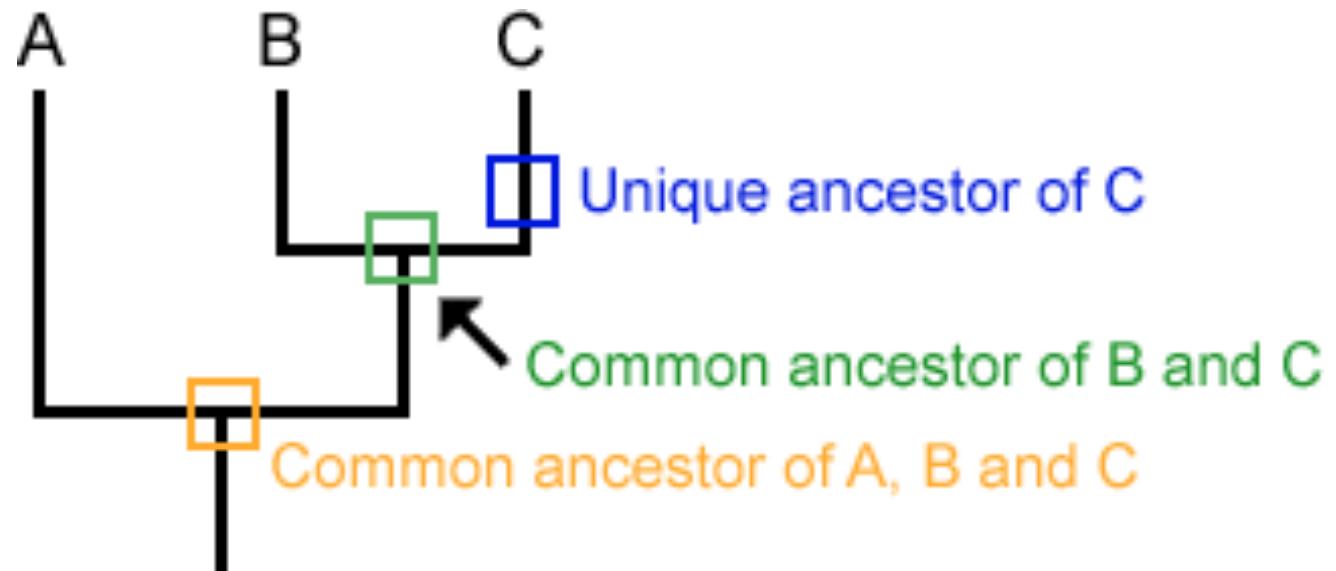
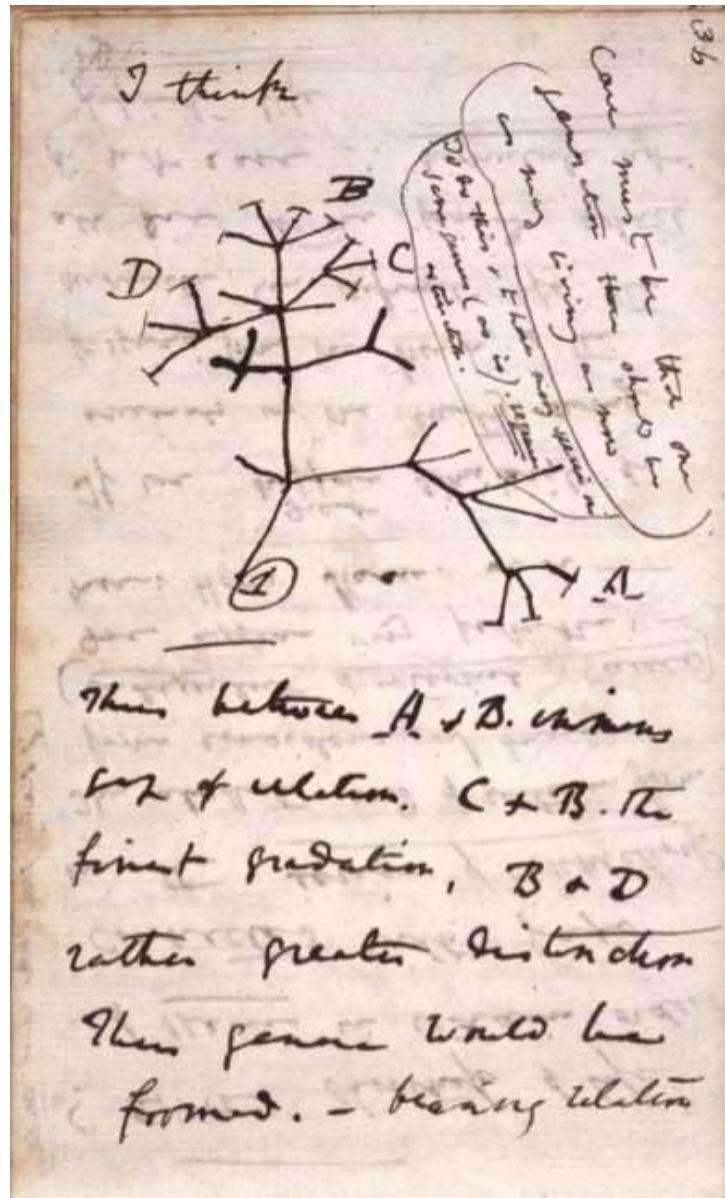


Evolutionary divergence and convergence in proteins
Zuckerkandl, E. and Pauling, L (1965)

A brief history of bioinformatics

Jeff Gauthier, Antony T Vincent, Steve J Charette, Nicolas Derome
Briefings in Bioinformatics (2018) <https://doi.org/10.1093/bib/bby063>

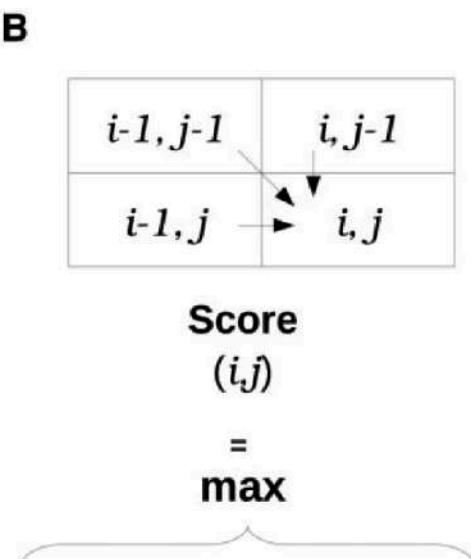
Relationships between sequences recapitulate evolutionary relationships



A mathematical framework for sequence alignments

A match +5 mismatch -4 gap -1

	A	T	C	G	
0	0	0	0	0	
A	0	5	-1	-1	-1
T	0	4	10 ← 9	8	
G	0	3	9	8	14



- Score $(i-1, j-1)$
+ Match / Mismatch
- Score $(i, j-1)$ + gap
- Score $(i-1, j)$ + gap

C
Best Alignment :
ATCG
|| |
AT G
(Score = 38)

Table 1. An excerpt of the PAM1 amino acid substitution matrix

10 ⁴ P ^a		Ala	Arg	Asn	Asp	Cys	Gln	...	Val
		A	R	N	D	C	Q	...	V
Ala	A	9867	2	9	10	3	8	...	18
Arg	R	1	9913	1	0	1	10	...	1
Asn	N	4	1	9822	36	0	4	...	1
Asp	D	6	0	42	9859	0	6	...	1
Cys	C	1	1	0	0	9973	0	...	2
Gln	Q	3	9	4	5	0	9876	...	1
...
Val	V	13	2	1	1	3	2	...	9901

^aEach numeric value represents the probability that an amino acid from the i-th column be substituted by an amino acid in the j-th row (multiplied by 10 000).

A brief history of bioinformatics

Jeff Gauthier, Antony T Vincent, Steve J Charette, Nicolas Derome
Briefings in Bioinformatics (2018) <https://doi.org/10.1093/bib/bby063>

1970-2000s – Paradigm shifts and parallel advances in biology and computer science

- Protein sequencing to DNA sequencing (faster / cheaper)
- Use DNA sequences to infer phylogenetic trees
- Sequence of marker genes and genomes
- Beyond sequences (structural bioinformatics)

- Faster computers
- GPUs
- Free software movement
- New Programming languages (Perl created by Larry Wall in 1987)

- Internet
- Online databases (NCBIs)

A brief history of bioinformatics

Jeff Gauthier, Antony T Vincent, Steve J Charette, Nicolas Derome
Briefings in Bioinformatics (2018) <https://doi.org/10.1093/bib/bby063>

Different sequencing platforms /
History of sequencing

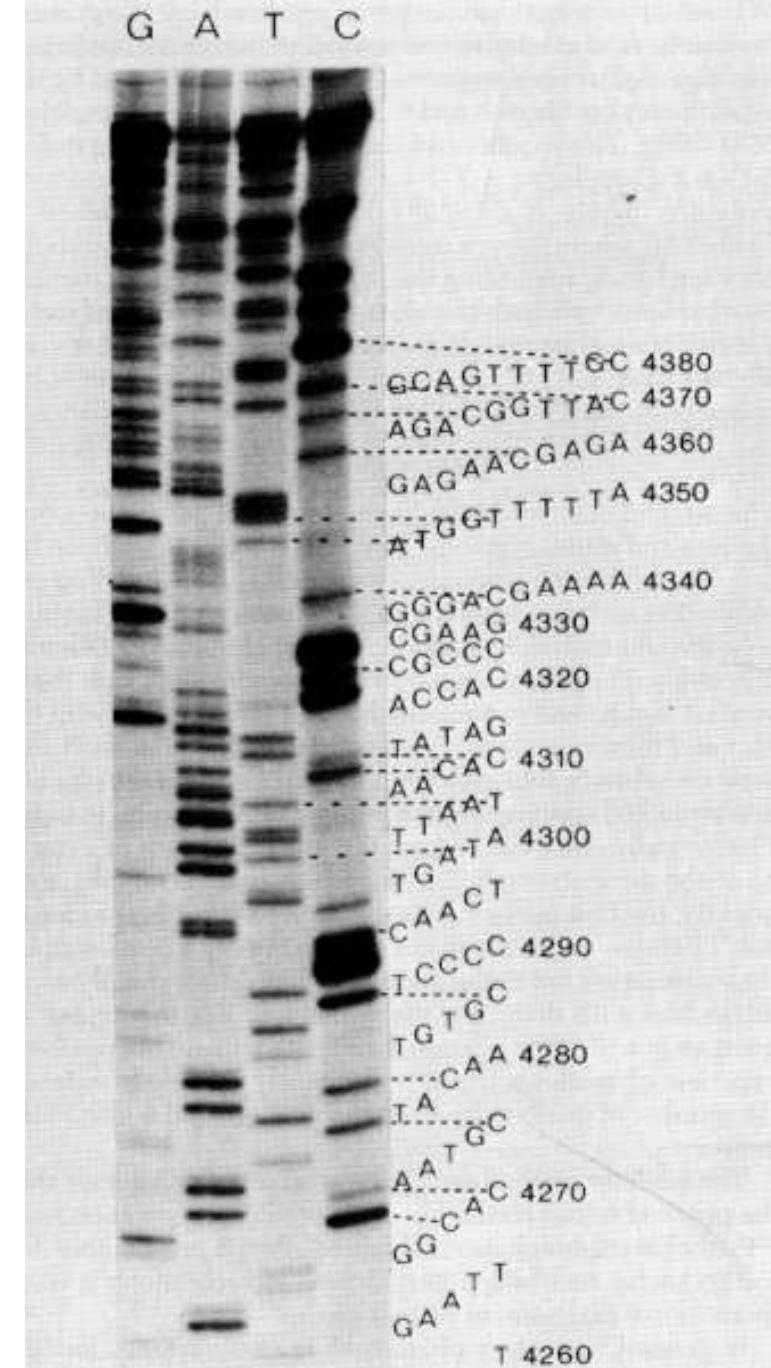
DNA sequencing with chain-terminating inhibitors

(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

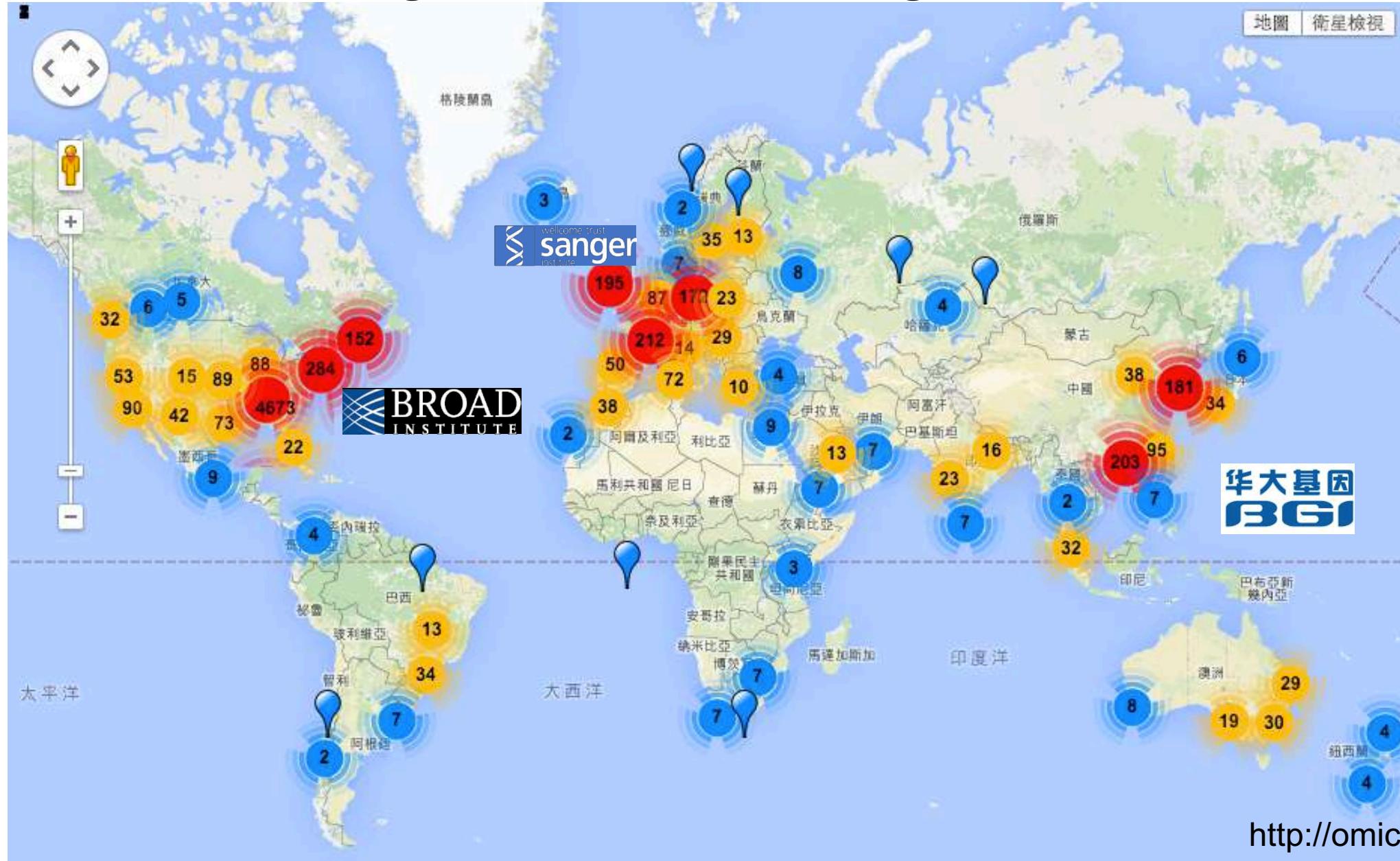
Contributed by F. Sanger, October 3, 1977



ABI 3730xi at TIGR (1.6Mb per day)



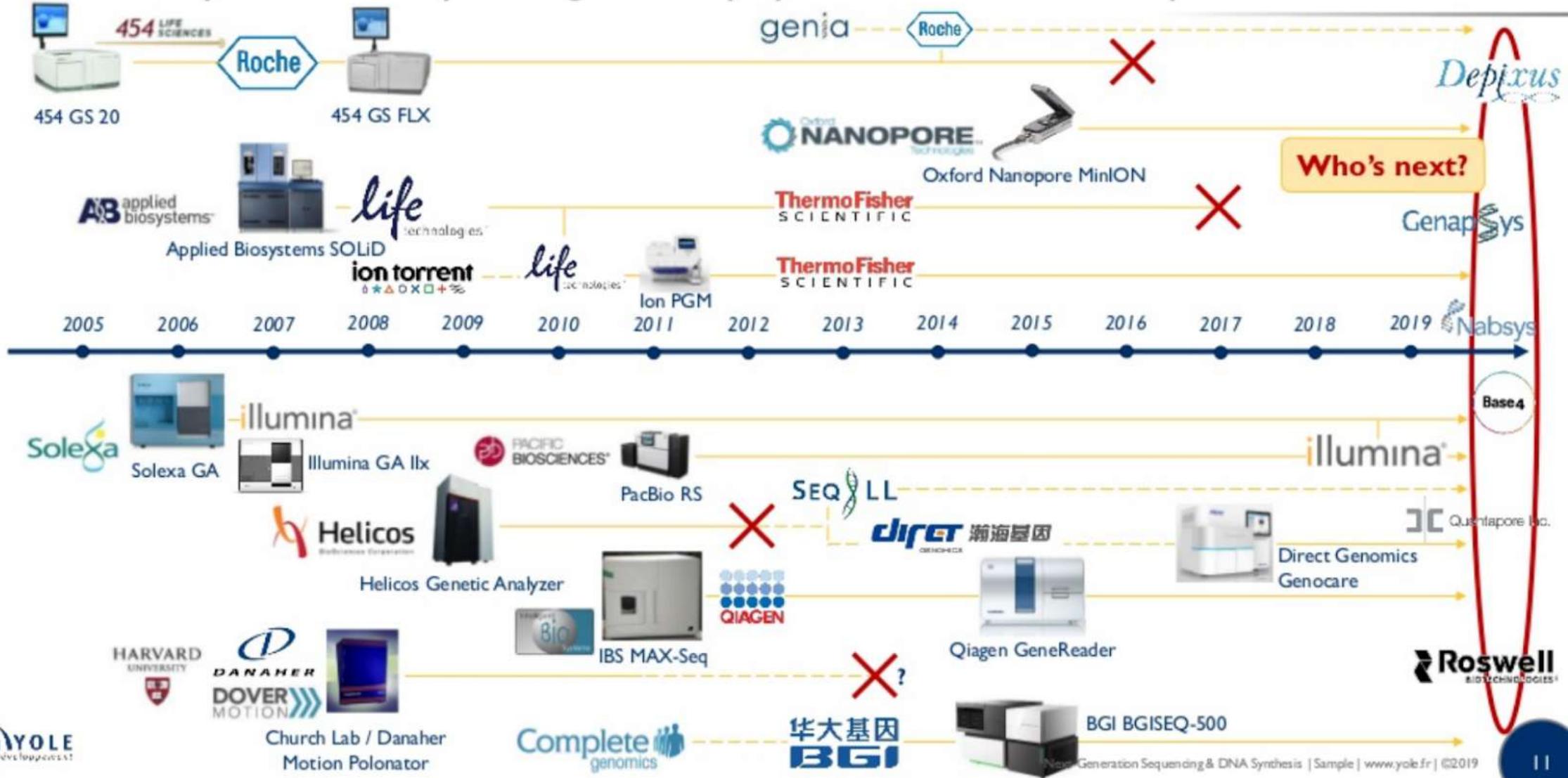
World competing for sequencing power



INTRODUCTION

Clip slide

History of DNA sequencing – Main players' first commercial products and M&A

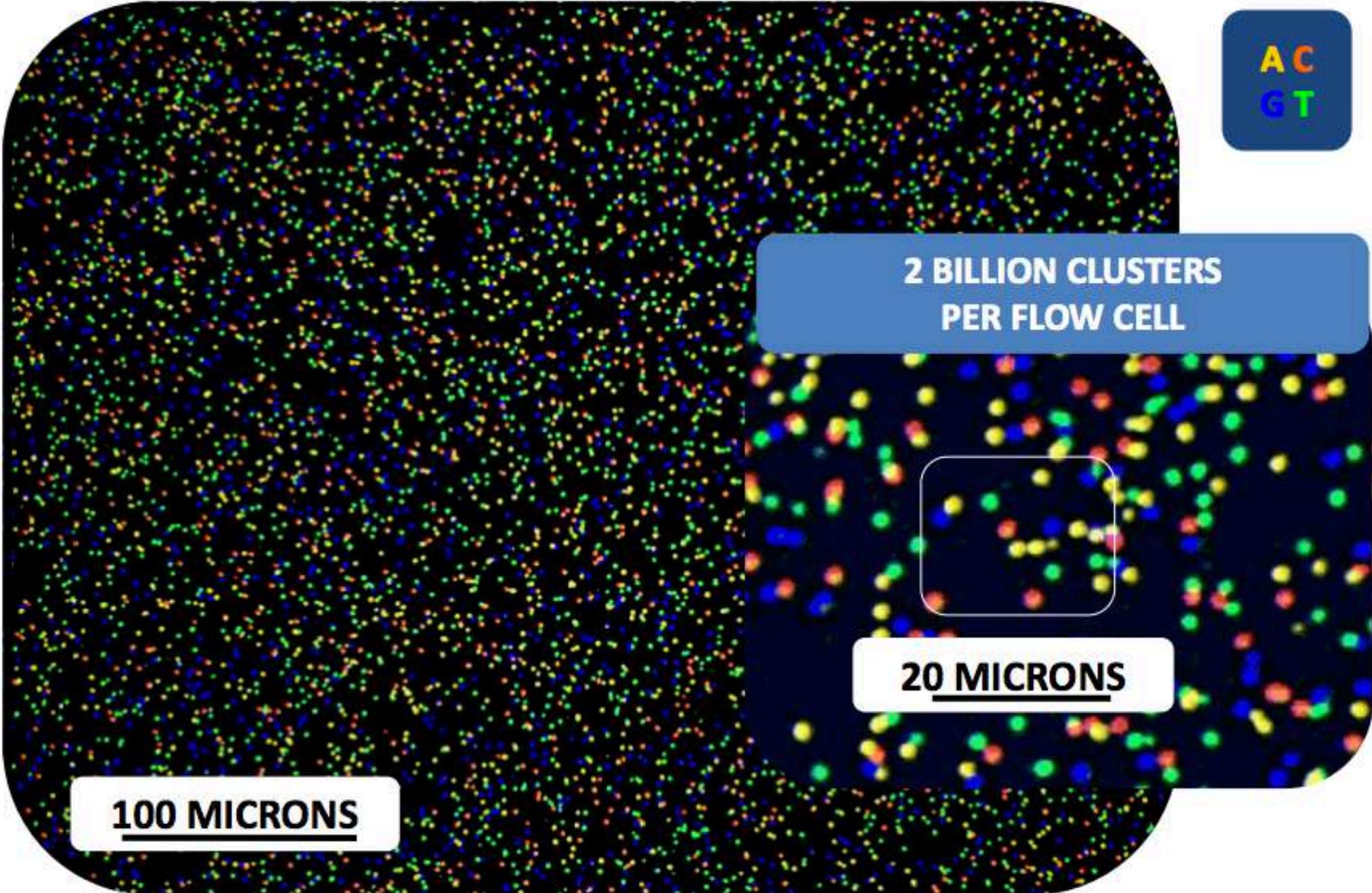


Illumina HiSeq



Sequencing by synthesis

<https://www.youtube.com/watch?v=fCd6B5HRaZ8>

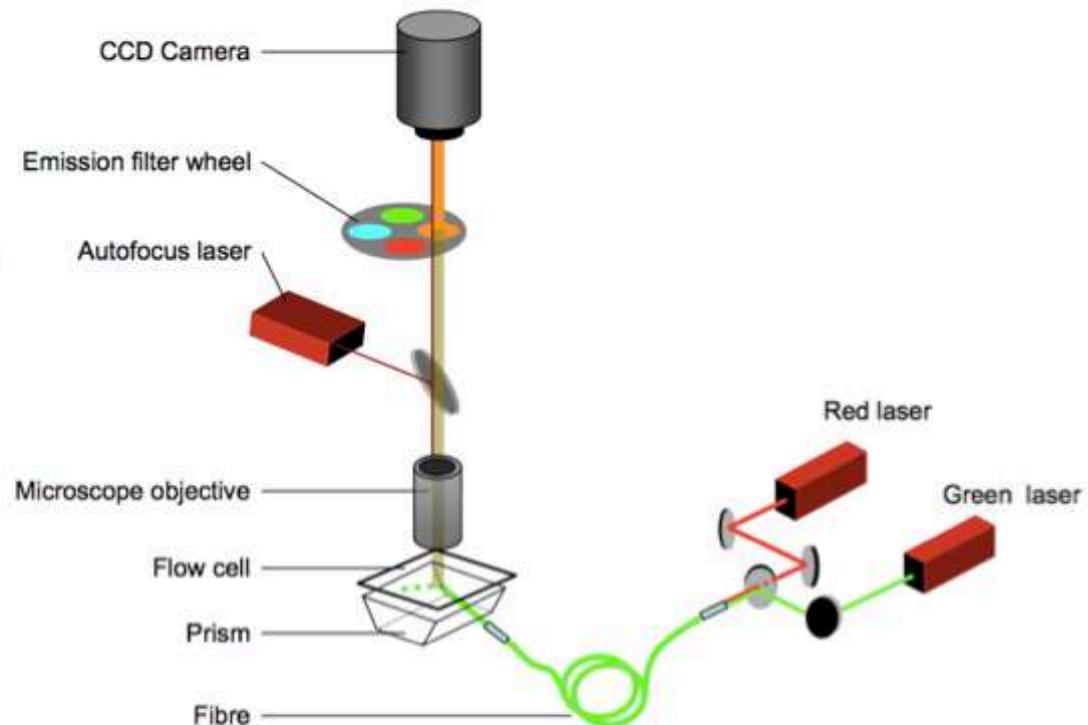
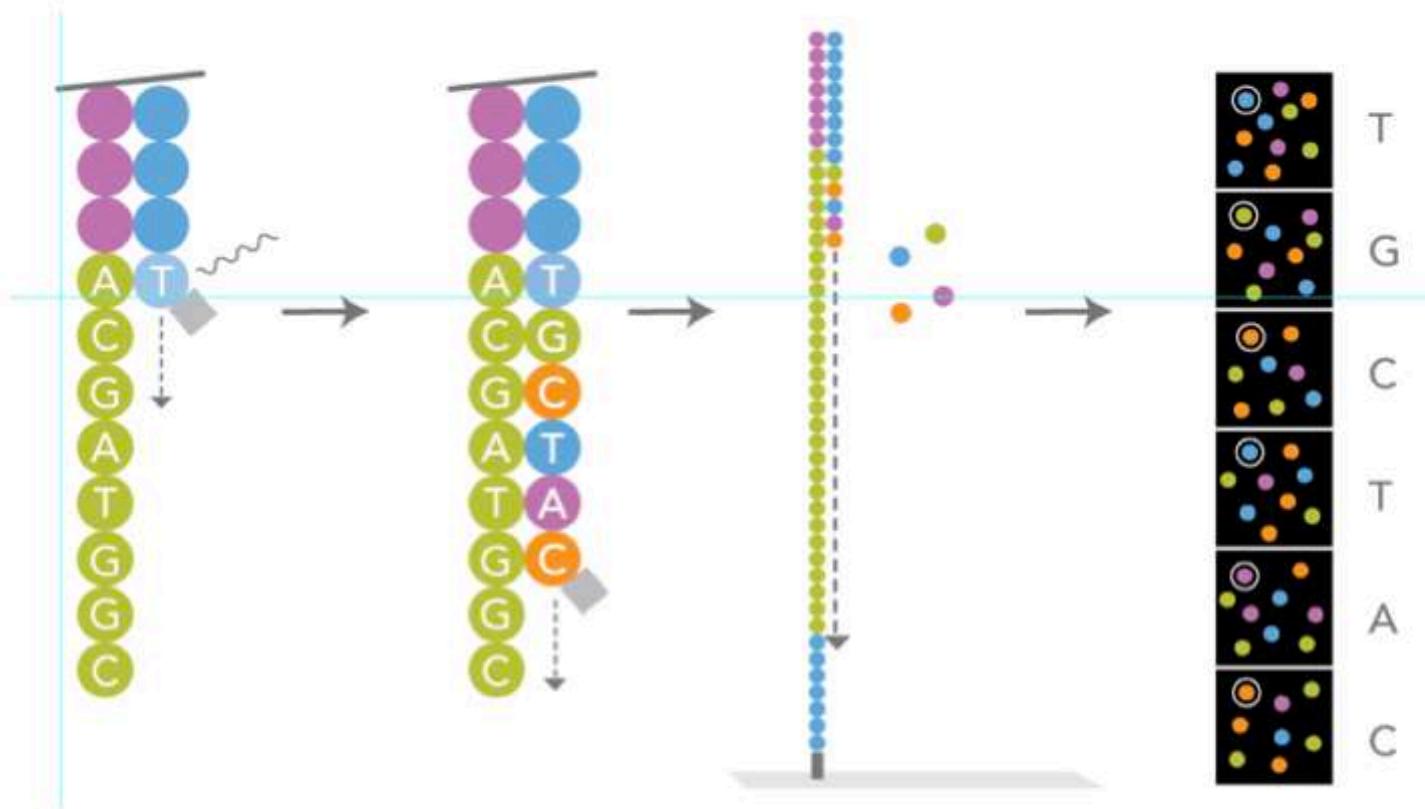


AC
GT

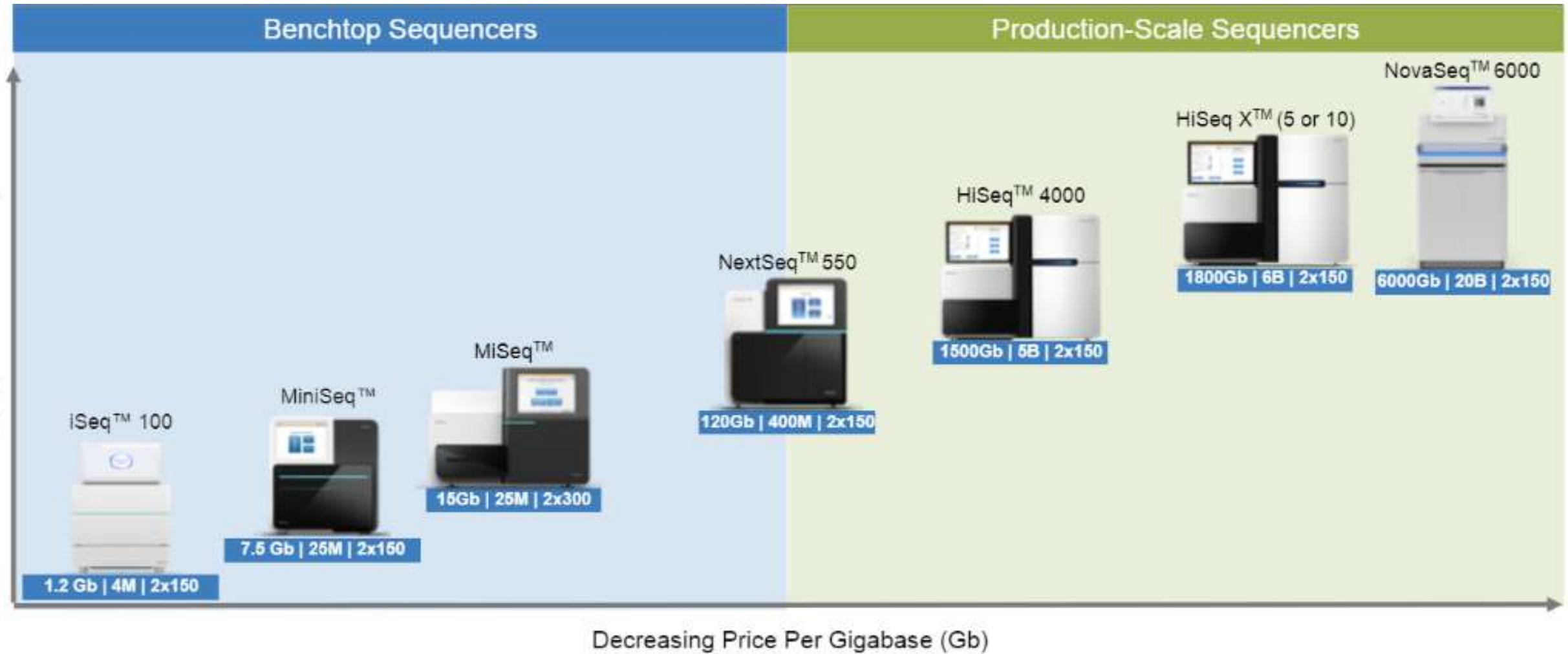
100 MICRONS

**2 BILLION CLUSTERS
PER FLOW CELL**

20 MICRONS



Illumina machines



And the arrival of 3rd generation sequencing...
(much longer read lengths and not so bad yield!!)

PacBio (Pacific Biosciences)



RSII



Sequel II

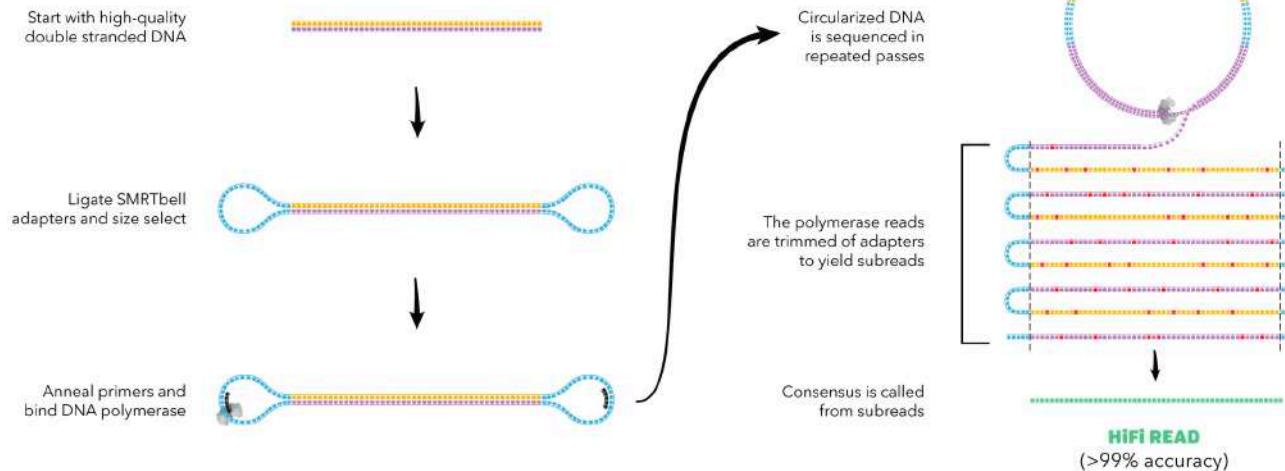
Single molecule sequencing

<https://www.youtube.com/watch?v=NHCJ8PtYCFc>

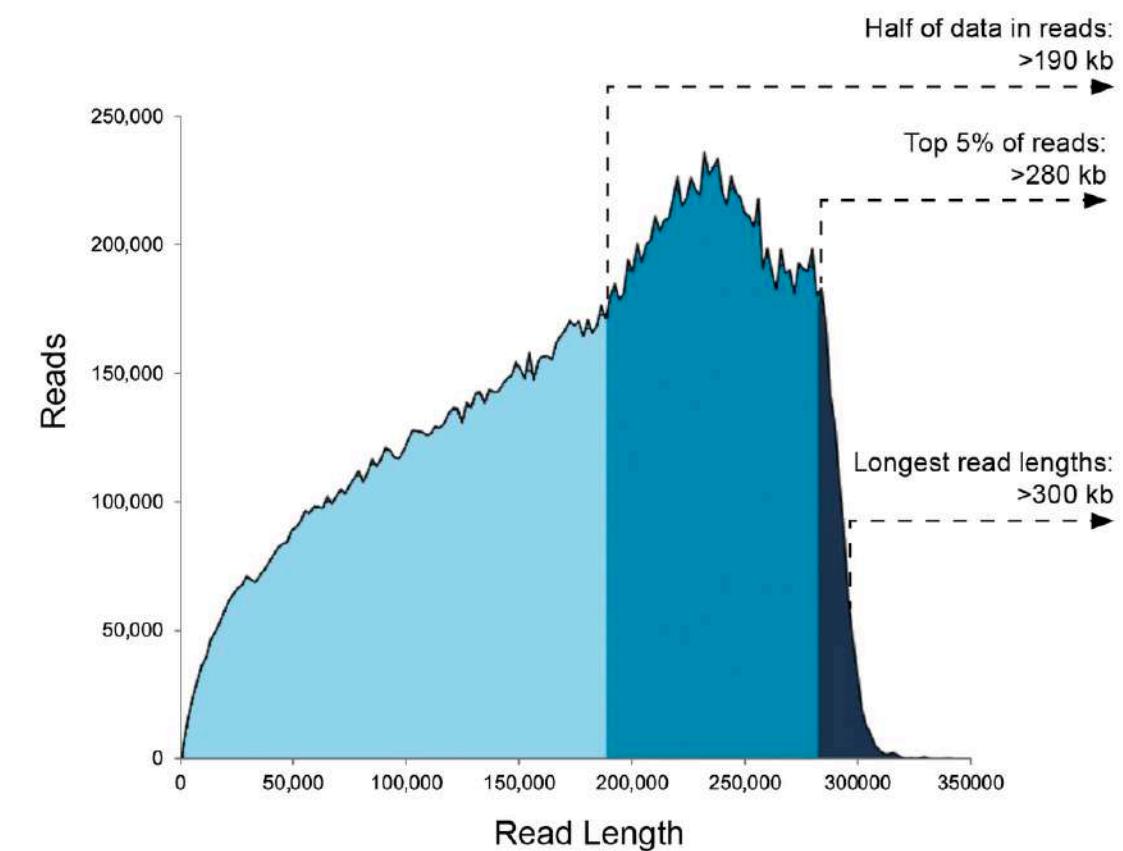
PacBio (Pacific Biosciences)



Produce HiFi reads using the circular consensus sequencing (CCS) mode to provide base-level resolution with >99% single-molecule read accuracy for the detection of all variant types from single nucleotide to structural variants. Learn more about the advantages of [long reads with high accuracy](#).



Half of data in reads: >190 kb
Data per SMRT Cell: Up to 50 Gb



Oxford Nanopore



Key	SmidgION	Flongle	MinION	GridION	PromethION
System Price	TBC	Included in \$5K Starter Pack	Included in \$1K Starter Pack	Included in \$50K Starter Pack	Included in \$135K Starter Pack
Number of channels	200 channels	128 channels	512 channels	$5 \times 512 = 2,560^*$	$48 \times 3,000^* = 144,000$
Per flow cell Current Data – Max Data	TBC	1 - 3.3 Gb	17 - 40 Gb	17 - 40 Gb	125 - 311 Gb
Per Device Current Data – Max Data				85 - 200 Gb	3/6 - 20 Tb
Price per Gb Current Data – Max Data	TBC	\$90 - \$30	\$30 - \$12.5	\$17.5 - \$7.5	\$5 - \$2

Oxford Nanopore – how it works

Introduction to nanopore

<https://vimeo.com/297106166>

Voltrax

<https://vimeo.com/297106291>

Sequencing for farmers

<https://vimeo.com/294216876>

@ Oceans

<https://vimeo.com/294744892>

Reference

<https://nanoporetech.com/how-it-works>

Nanopore Sequencing of Ebola Viruses Under Outbreak Conditions

<https://www.youtube.com/watch?v=SYBzPEoENWI> ; <https://www.nature.com/articles/nature16996>

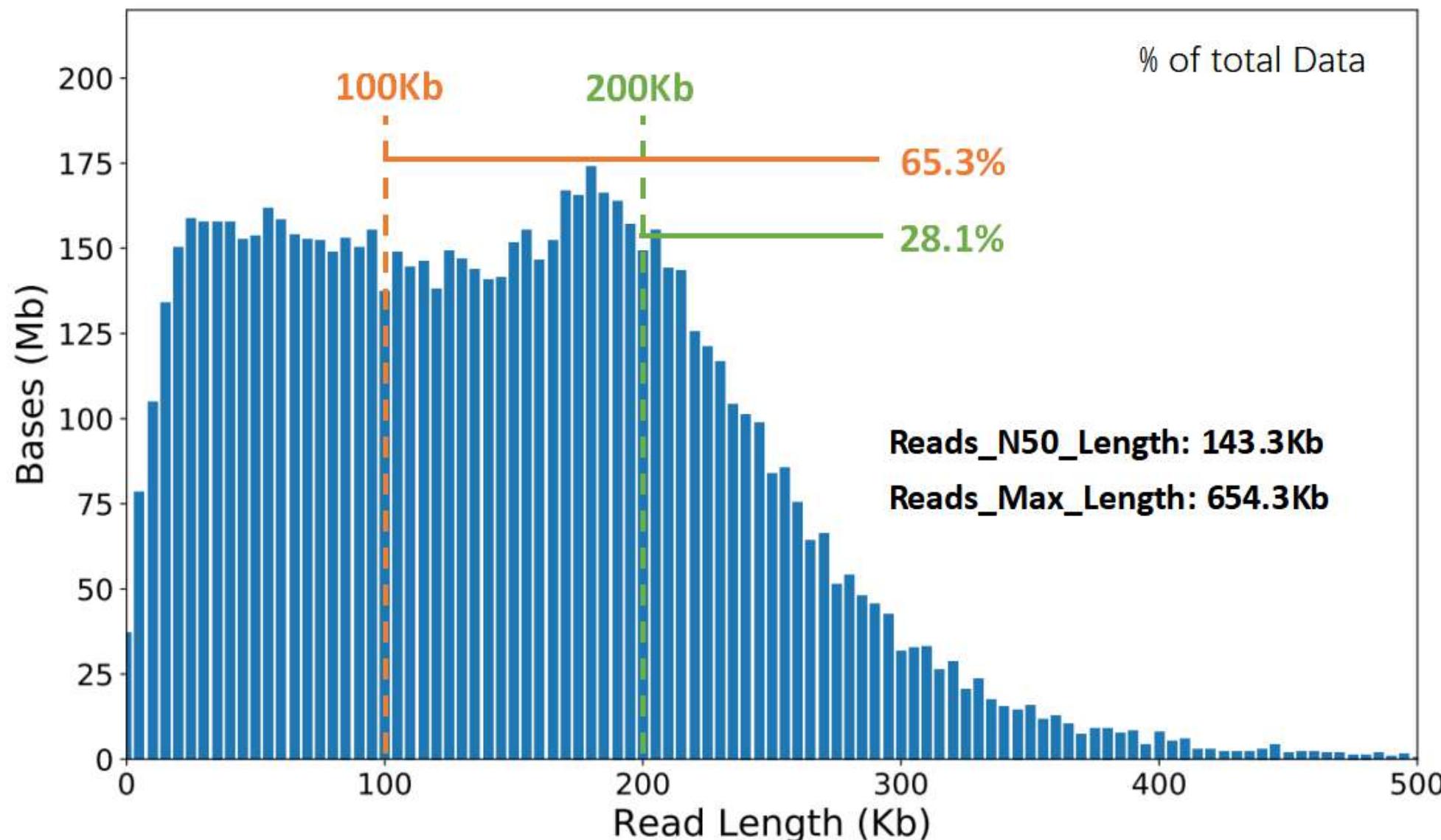
Rainforest

<https://www.youtube.com/watch?v=6RRSxWtJPUw>

From Extreme to everyday

https://www.youtube.com/watch?v=tQ_oo7_36r8

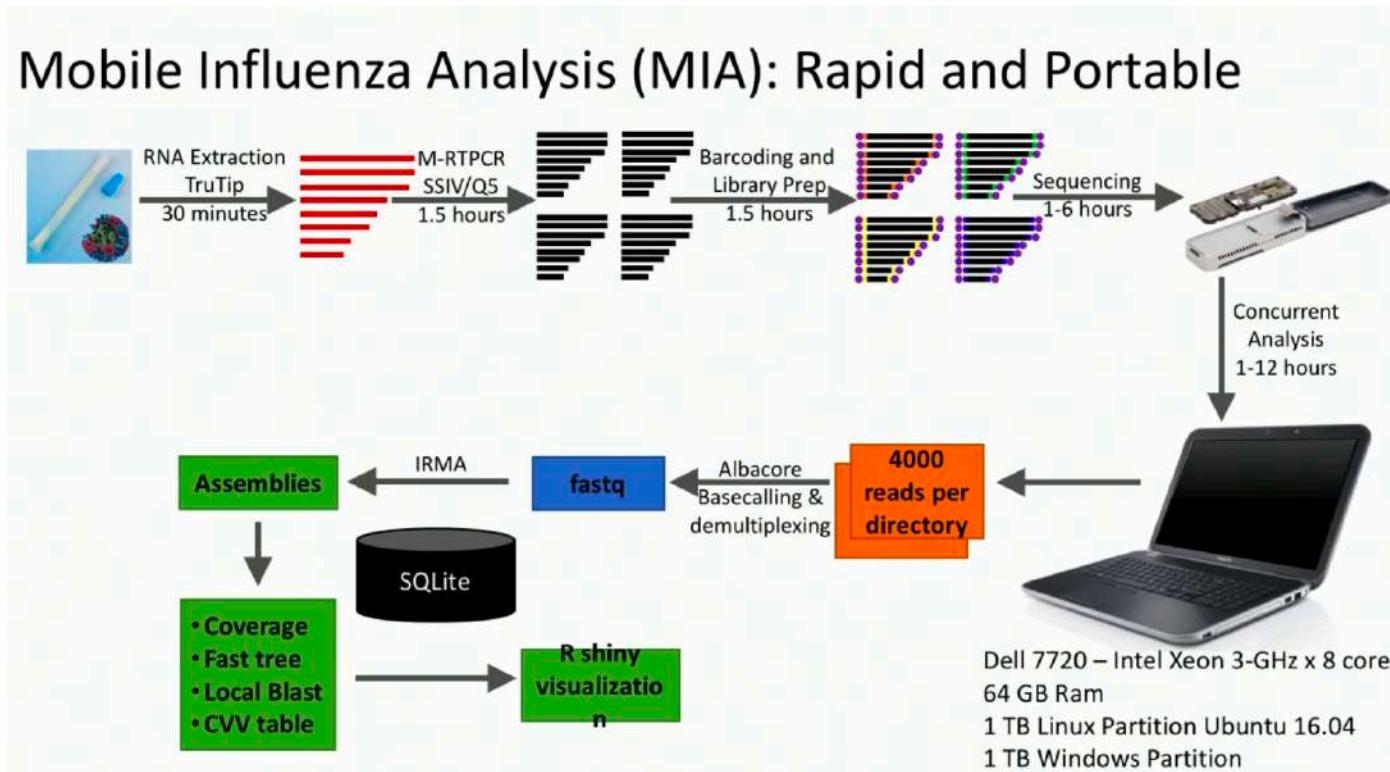
Read length and capacity go beyond



Mobility of sequencing



Matthew Keller: Deployable NGS for Influenza virus field surveillance and outbreak response



<https://nanoporetech.com/resource-centre/matthew-keller-deployable/ngs-influenza-virus-field-surveillance-and-outbreak>

Matthew Keller: Deployable NGS for Influenza virus field surveillance and outbreak response

Tom Connor (@tomrconnor)

Within 24 hours of starting the @NetworkArtic PCR protocol on our viral extracts for @nanopore sequencing our first sequences were up, available on GISAID and had already been analysed by Nextstrain. Incredible. @WalesMicrobiol @SmallRedOne @GenomicsWales [twitter.com/nextstrain/sta...](https://twitter.com/nextstrain/status/1234567890)

Nextstrain @nextstrain
Replying to @nextstrain

The two genomes from Wales each group the large European outbreak clade, but don't group together, suggesting separate introductions. Thanks to @SmallRedOne, @tomrconnor, @PublicHealthW, @WalesMicrobiol 2/3

Phylogeny

Adaptive evolution

NETT LINEUP

Wales/PWAT/1/2020 ←
Brazil/MSBR/61/2020
Niger/Legon/01/2020
Kenya/LGIG/2020
Venezuela/DOHE/01/2020
Germany/Bavaria/Wittenberg/1/2020
Germany/Bavaria/Wittenberg/1/2020
Venezuela/DOHE/01/2020
Venezuela/DOHE/01/2020
Spain/ICV/001/2020
Wales/PWAT/2/2020 ←
Switzerland/2000/77/37/2020
Italy/UMS/3/2020
Germany/Bavaria/3/2020

Geography

NETT 2020

Map of Europe showing the location of Wales in the United Kingdom. A yellow dot marks the location of Wales on the map.

SARS-CoV-2 Whole genome sequencing



7hr

RNA to answer

Of which ~1 hr
sequencing time

<https://nanoporetech.com/about-us/news/novel-coronavirus-covid-19-information-and-updates>

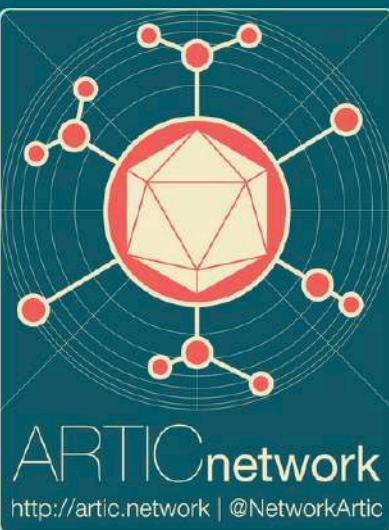
<https://nanoporetech.com/about-us/news/covid19-community>

About

The Project

This project is developing an end-to-end system for processing samples from viral outbreaks to generate real-time epidemiological information that is interpretable and actionable by public health bodies. Fast evolving RNA viruses (such as Ebola, MERS, SARS, influenza etc) continually accumulate changes in their genomes that can be used to reconstruct the epidemiological processes that drive the epidemic. Based around a recently developed, single-molecule portable sequencing instrument, the Oxford Nanopore Technology MinION, we are creating a ‘lab-in-a-suitcase’ that can be deployed to remote and resource-limited locations. Targeting a wide-range of emerging viral diseases, the sequencing generation will be closely linked to the analysis platform to integrate these data and associated epidemiological knowledge to reveal the processes of transmission, virus evolution and epidemiological linkage with extremely rapid turn-around. This real-time approach will provide actionable epidemiological insights within days of samples being taken from patients.

nCoV-2019



There is a pressing need to understand more about the short-term genomic epidemiology and evolution of the recently described novel coronavirus (nCoV-2019). Initial cases were in Wuhan City, Hubei Province, China but now cases have been confirmed both more widely in China and internationally.

Viral genome data generated prospectively during outbreaks can help provide information about relatedness to other viruses, mode and tempo of evolution, geographical spread and adaptation to human hosts. This information can be used to assist in epidemiological investigations, particularly when combined with other types of data (e.g. case counts).

The ARTIC network is making available a set of materials (see below) to assist groups in sequencing the virus including a set of primers, laboratory protocols, bioinformatics tutorials and datasets. These are mainly focused around the use of the portable Oxford Nanopore MinION sequencer, although aspects of the protocol such as the primer scheme and sample amplification may be generalised to other sequencing platforms.



HELP DOCS BLOG LOGIN

Nextstrain

Real-time tracking of pathogen evolution

Nextstrain is an open-source project to harness the scientific and public health potential of pathogen genome data. We provide a continually-updated view of publicly available data alongside powerful analytic and visualization tools for use by the community. Our goal is to aid epidemiological understanding and improve outbreak response. If you have any questions, or simply want to say hi, please give us a shout at hello@nextstrain.org.



Scenarios now and then

1. [lab/hospital/mountain/sea] Collect samples (1.1, 1.2, 1.3...)
2. [lab/hospital] Extract DNA (2.1, 2.2, 2.3...)
3. [lab/hospital/company] Sequencing (3.1, 3.2, 3.3...)
4. [lab/company] Analysis
5. [lab/hospital] Report

Weeks

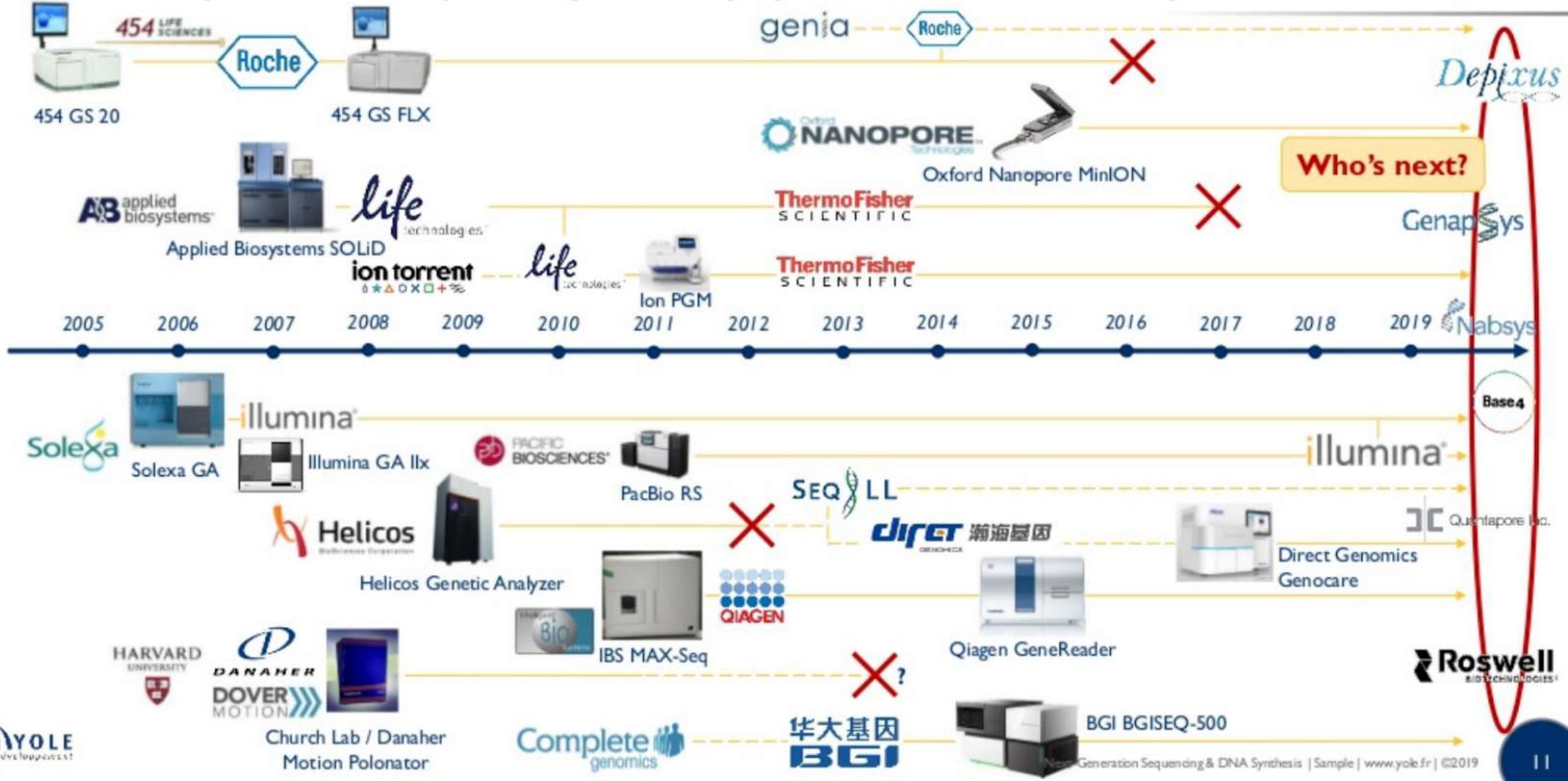
1. [lab/hospital/mountain/sea] Collect samples -> report

Minute

INTRODUCTION

Clip slide

History of DNA sequencing – Main players' first commercial products and M&A



Break here

Three 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

Counting:

- For a given region (gene) we want to know how much. → gene expression or metagenomics

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
	Insert size	The length of the DNA fragment
	Contig	A set of overlapping DNA segments that represents a consensus region of DNA
	Scaffold	Contigs separated by gaps of known length
	Coverage	The number of times a specific position in the genome is covered by reads

What is an alignment? (mapping)

Align the following two sequences:

ATTGAAAGCTA

GAAATGAAAAGG

1:

--ATTGAAA-GCTA

| | | | | |

GAAATGAAAAGG--

Scoring scheme is needed:

1 for match

-1 for mismatch

-2 for gap

2:

ATTGAAA-GCTA---

| | | | | |

---GAAATGAAAAGG

insertions / deletions (indels) mismatches

Which alignment is better?

Assembly

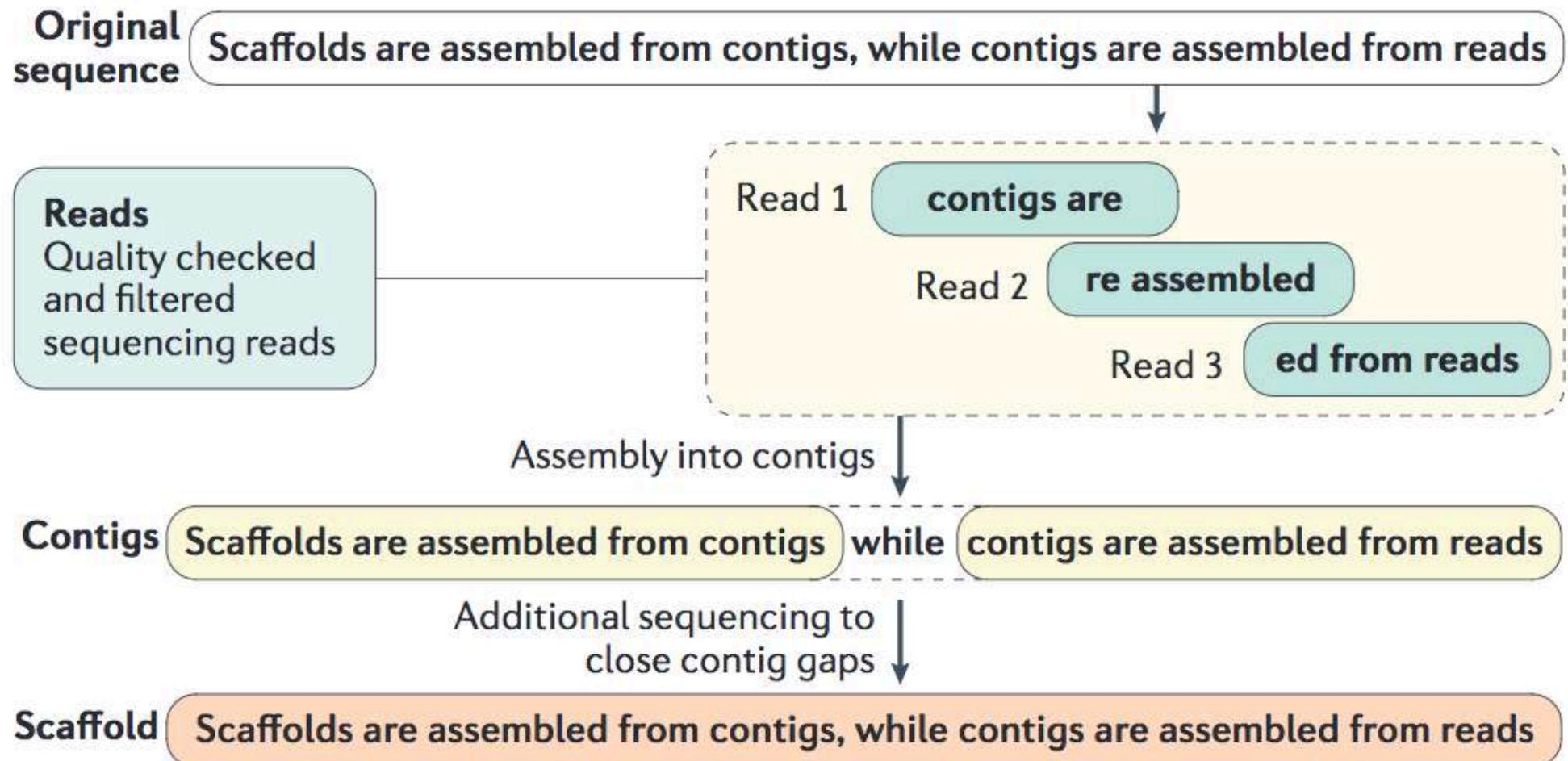
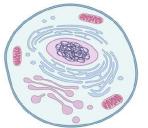


Figure 2 | **Sequence read assembly.** A mock example explaining bioinformatic sequence assembly along with the terms sequence, reads, contigs and scaffolds.

Assembly



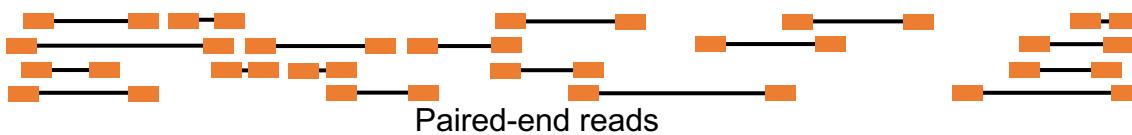
Genome



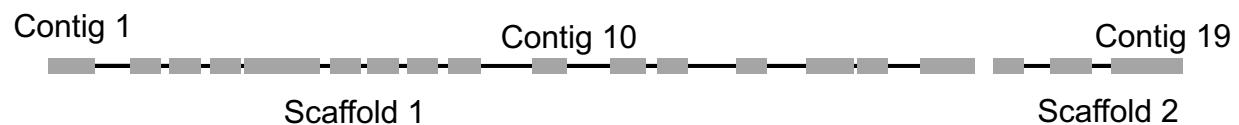
Fragment



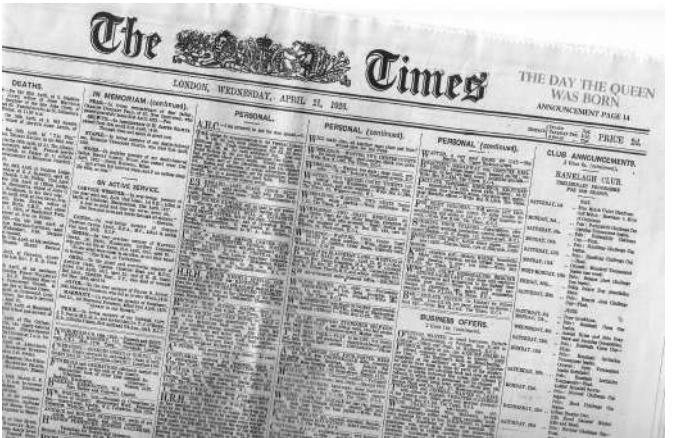
Sequence



Assemble



Assembly



Sequencing



Assembly



Genome
(3.000.000 letters)



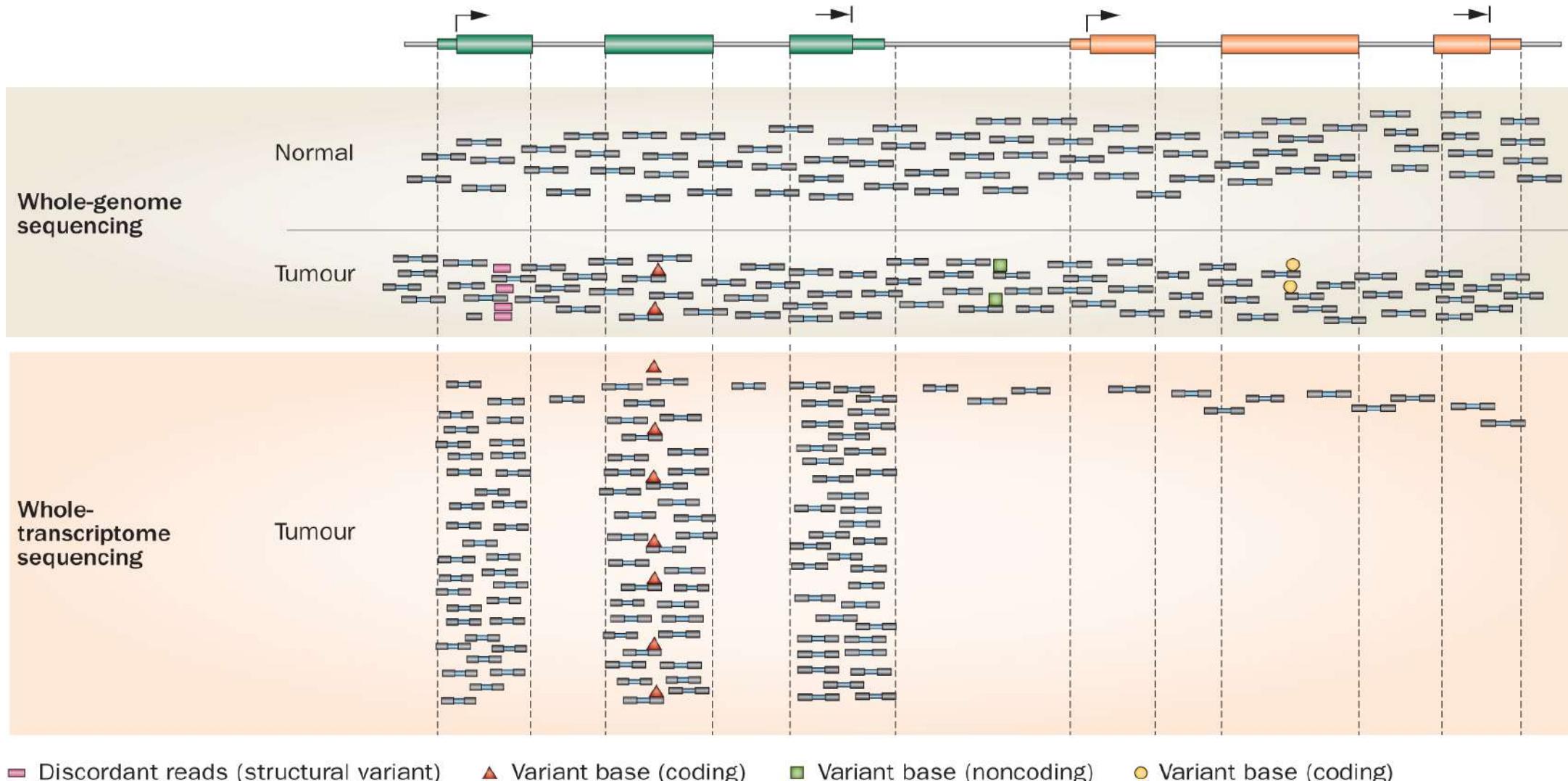
Genome
(3.000.000 letters)

Reads

(50-500 letters each)

Mapping

Reference genome depicting two example genes



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.

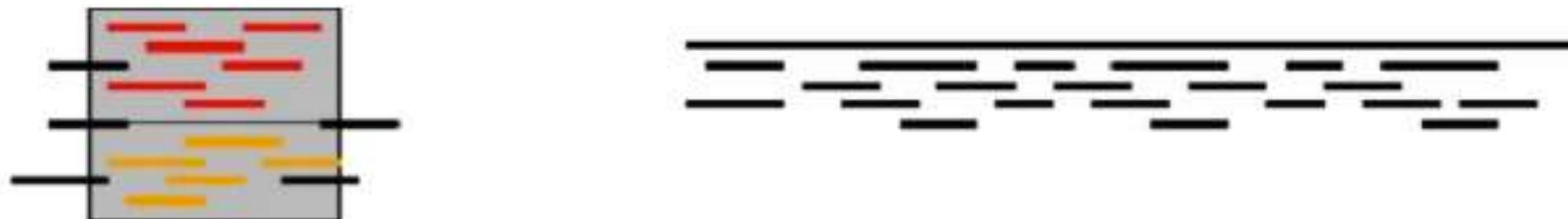
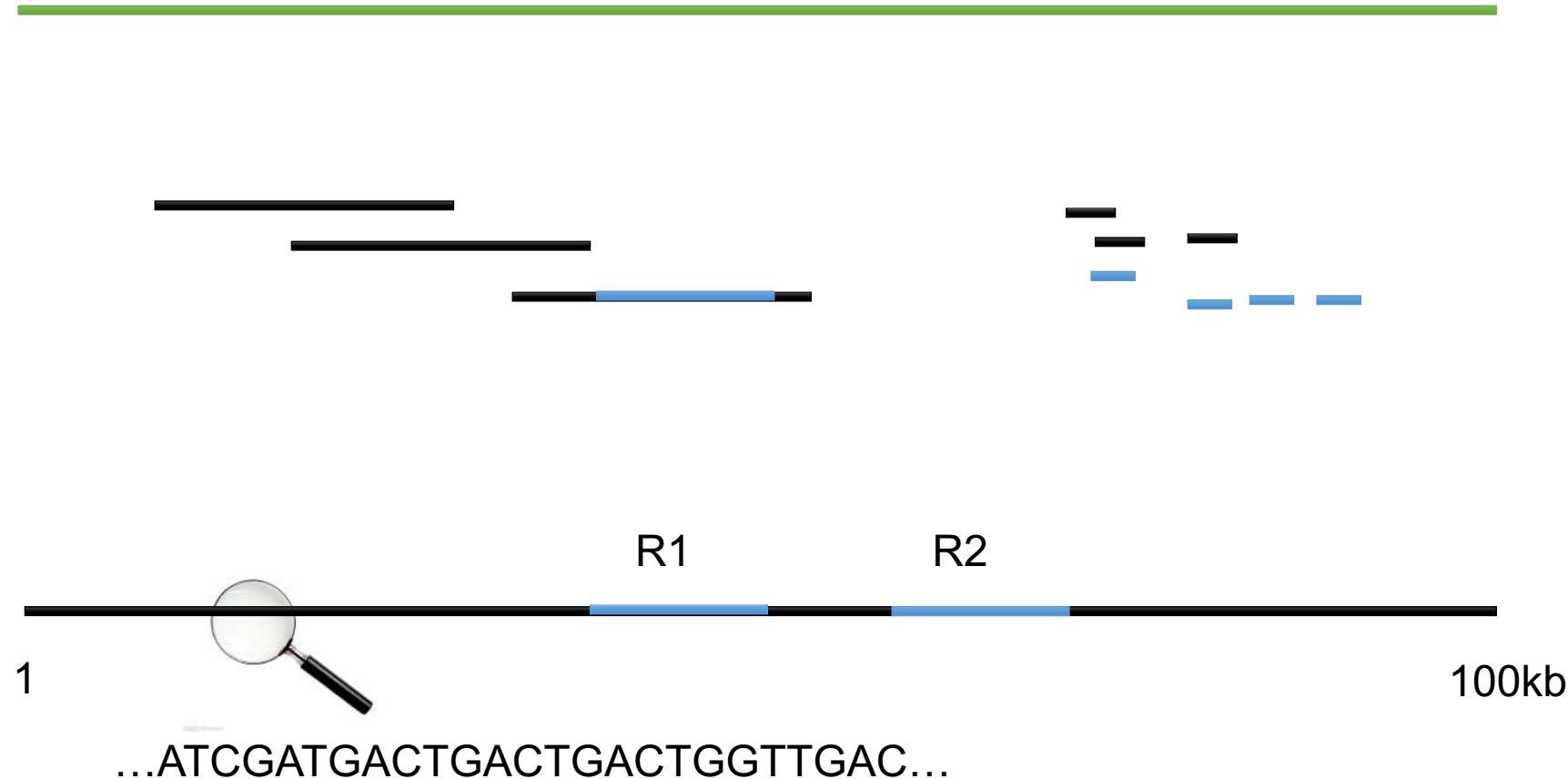
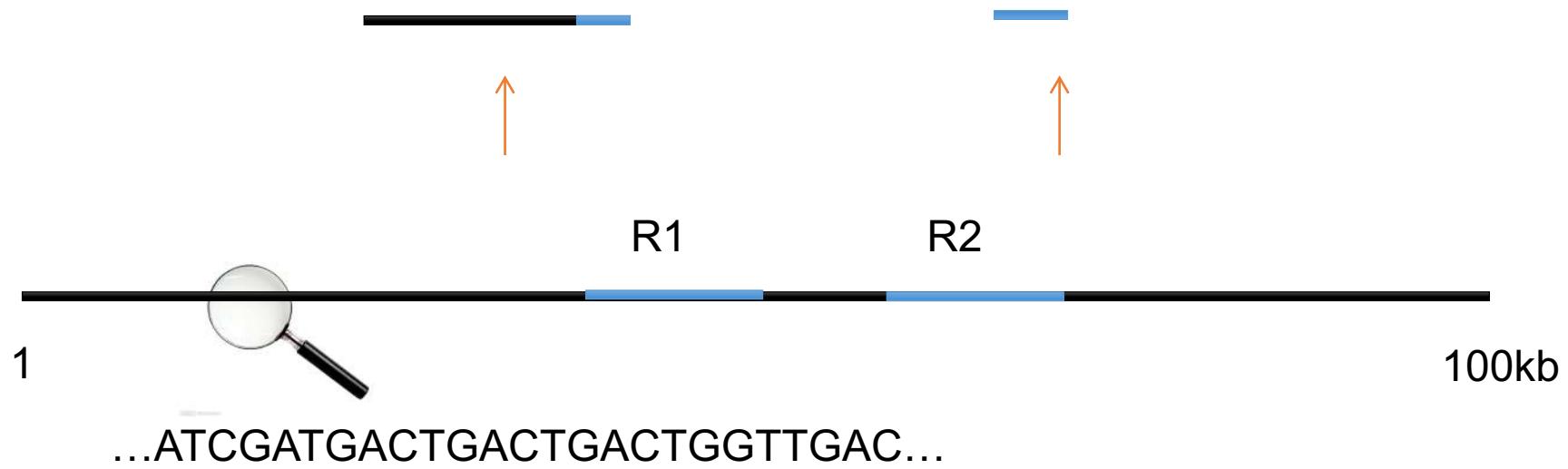


Figure 6. Genome mis-assembled due to a repeat. The assembly program incorrectly combined the reads from the two copies of the repeat leading to the creation of two separate contigs

Read length matters in sequencing



Paired end and insert size matter in sequencing



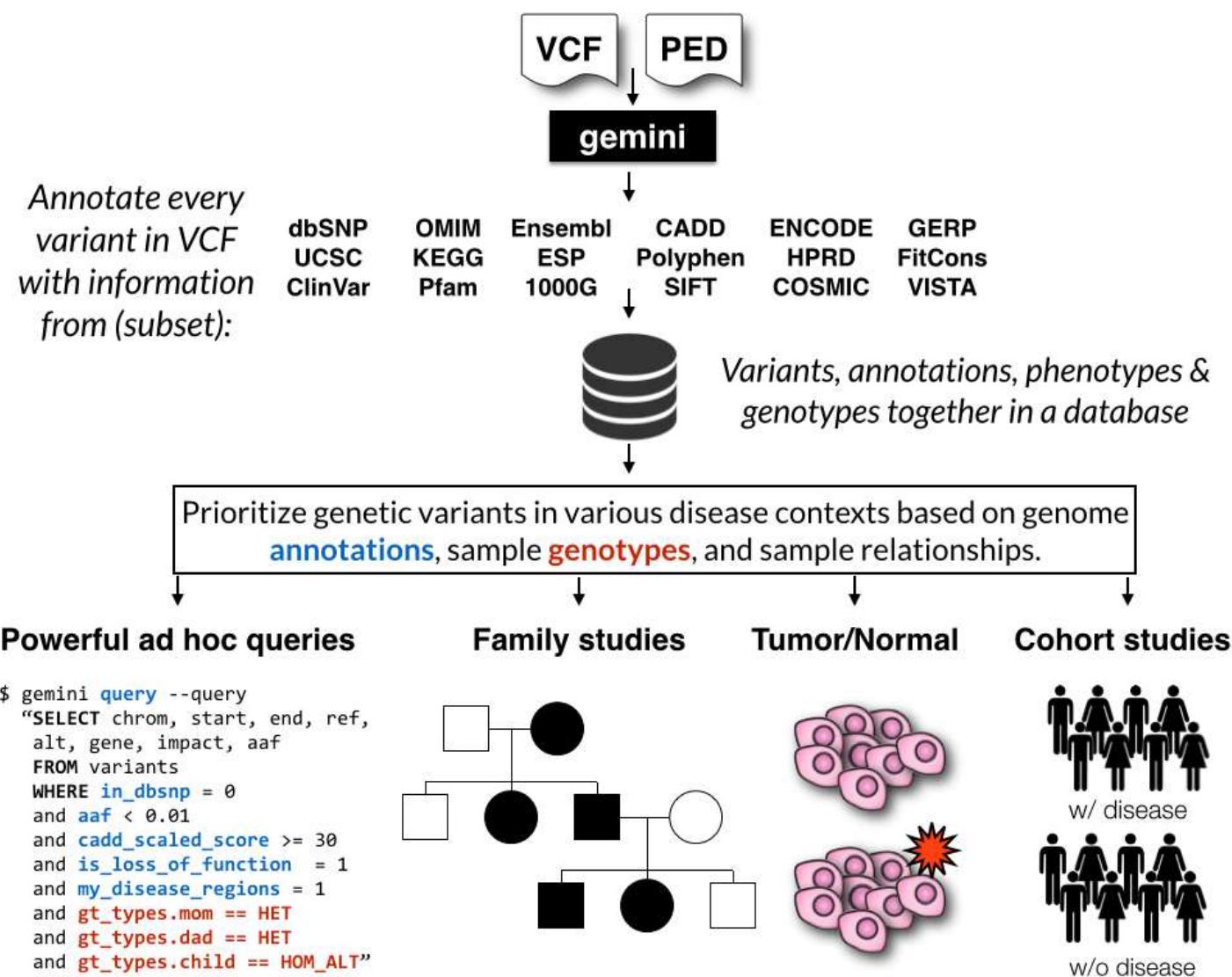
Depth matters in sequencing

	ATCGATGACTGACTGA ATGGTTGAC
	ATCGATGACTGACTGA ATGGTTGAC
	ATCC C ATGACTGACTGA ATGGTTGAC
	ATCGATGACTGACTGA ATGGTTGAC
	ATCGATGACTGACTGA ATGGTTGAC
	ATCGATGACTGAG TGA ATGGTTGAC
	ATCGATGACTGAG TGA ATGGTTGAC
	ATCGATGACTGAG TGA ATGGTTGAC
	ATCGATGACTGAG TGA ATGGTTGAC
10X	ATCGATGACTGAG TGA ATGGTTGAC
1X	Homozygous? Heterozygous? ATCGAT C ACTGACTGACTGGTTGAC

...ATCGATGACTGACTGACTGGTTGAC...

reference

Filtering and annotating variants



ClinVar
SIFT
Function
MAF (Minor Allele frequency)

Case studies

Classical genetics

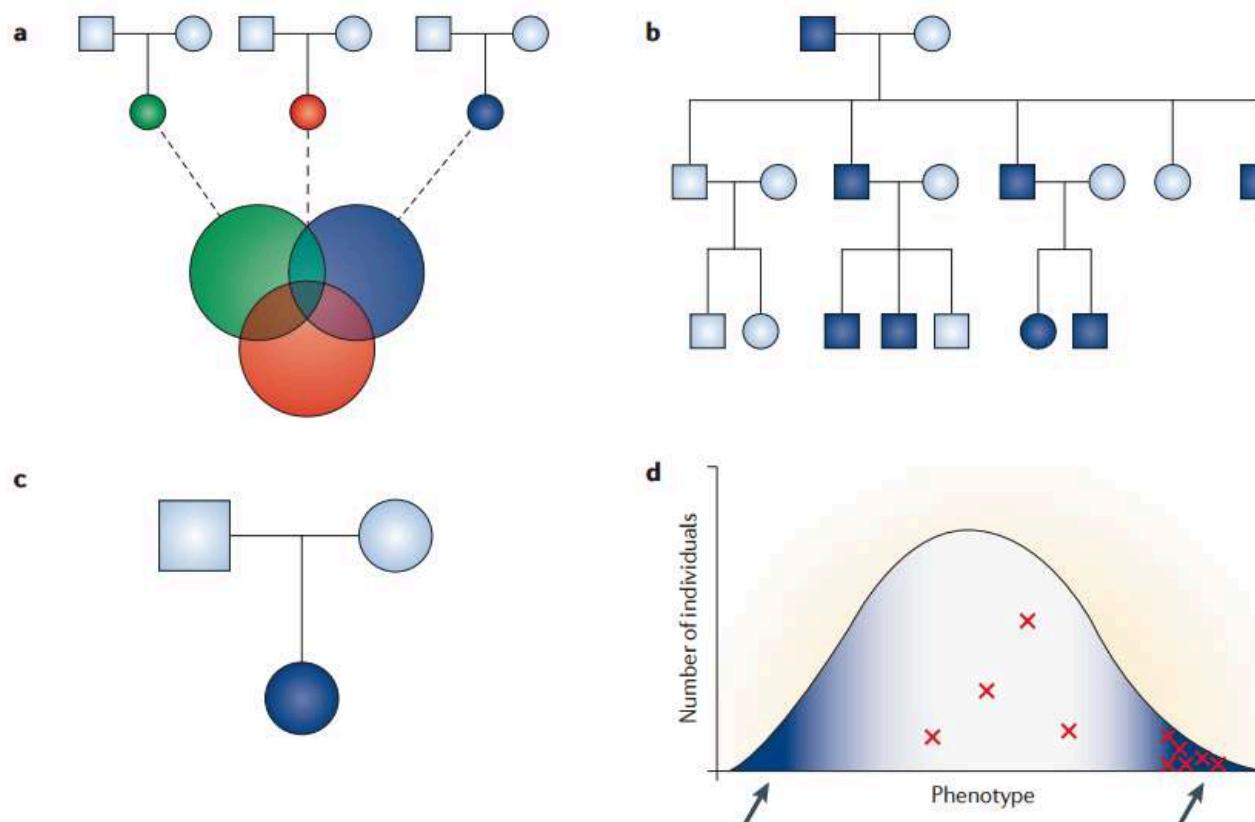
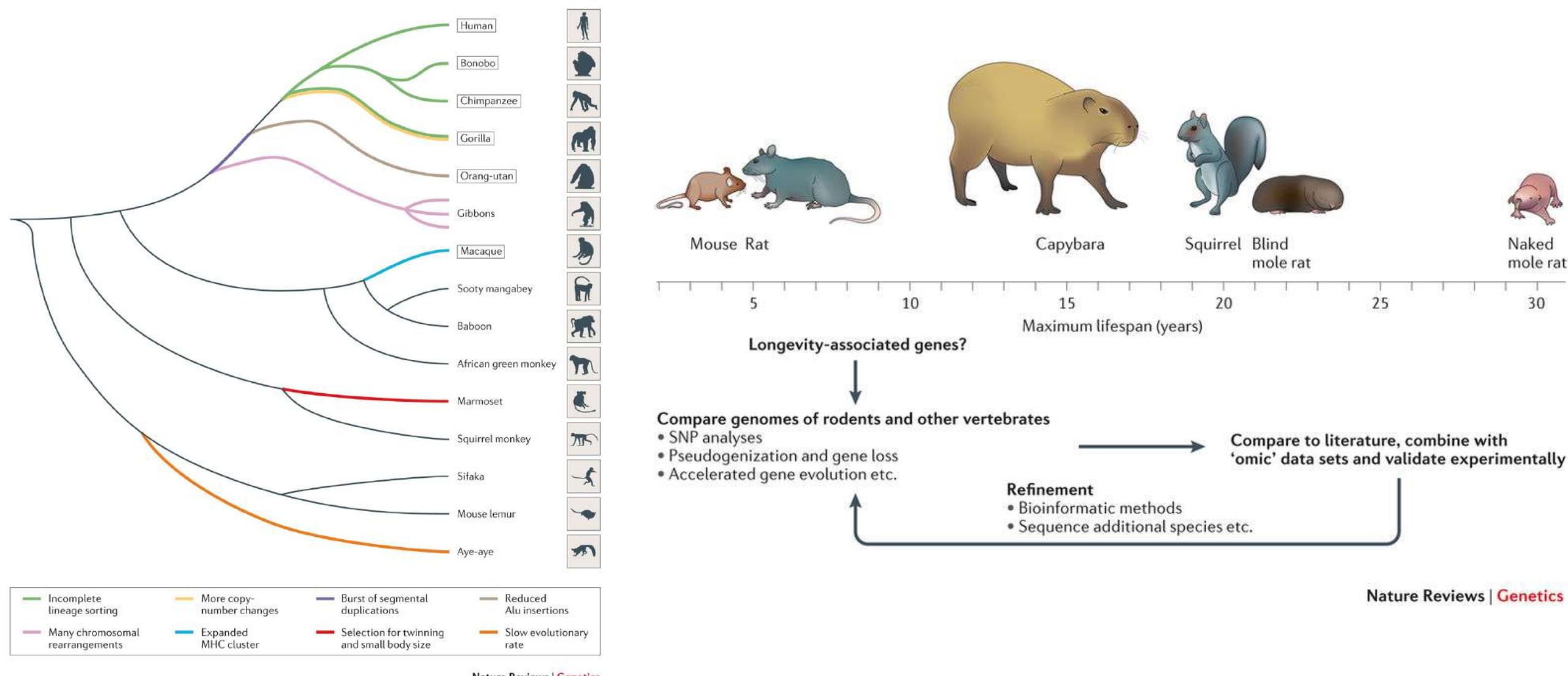


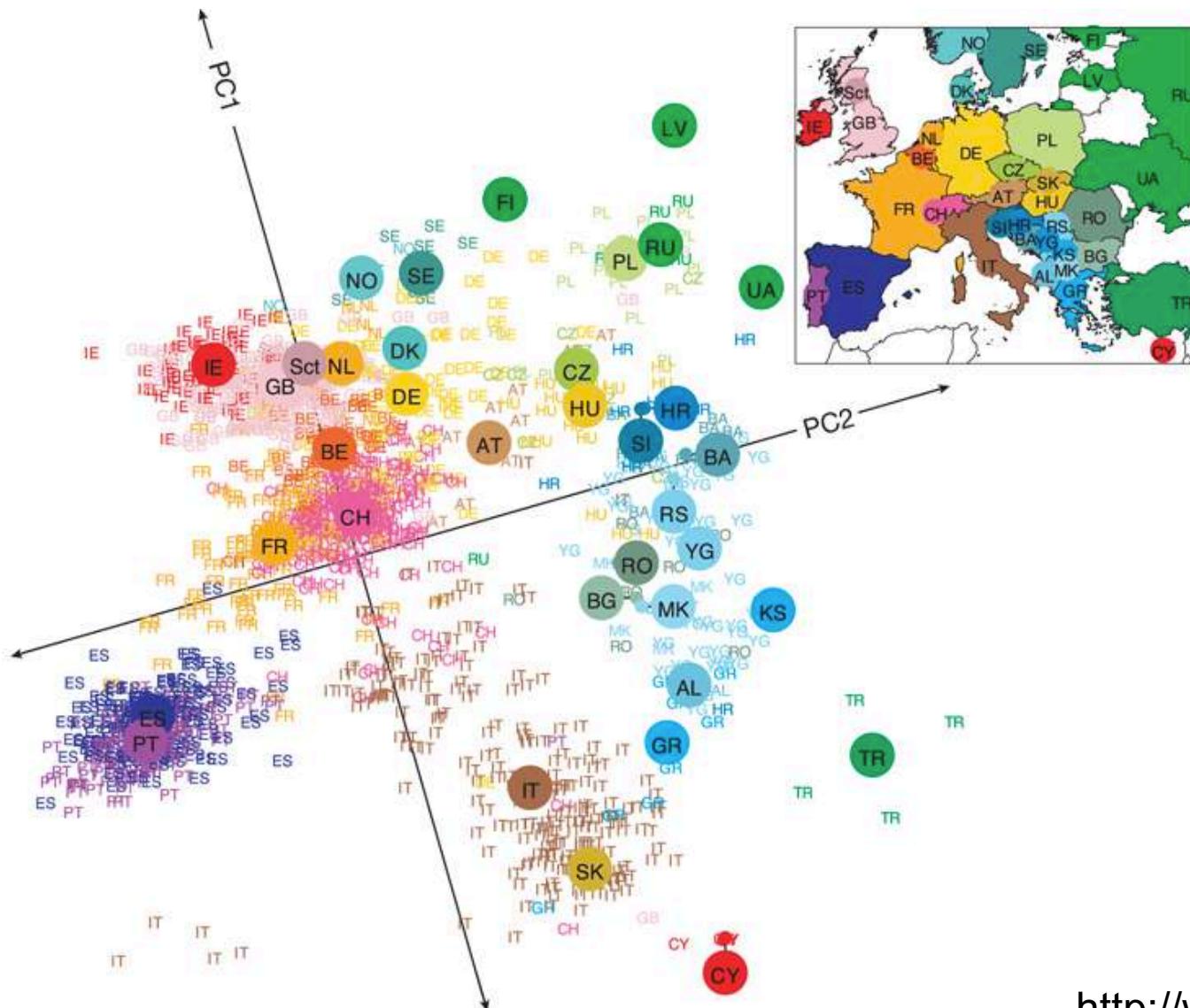
Figure 2 | Strategies for finding disease-causing rare variants using exome sequencing. Four main strategies are illustrated. **a** | Sequencing and filtering across multiple unrelated, affected individuals (indicated by the three coloured circles). This approach is used to identify novel variants in 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

Population genomics

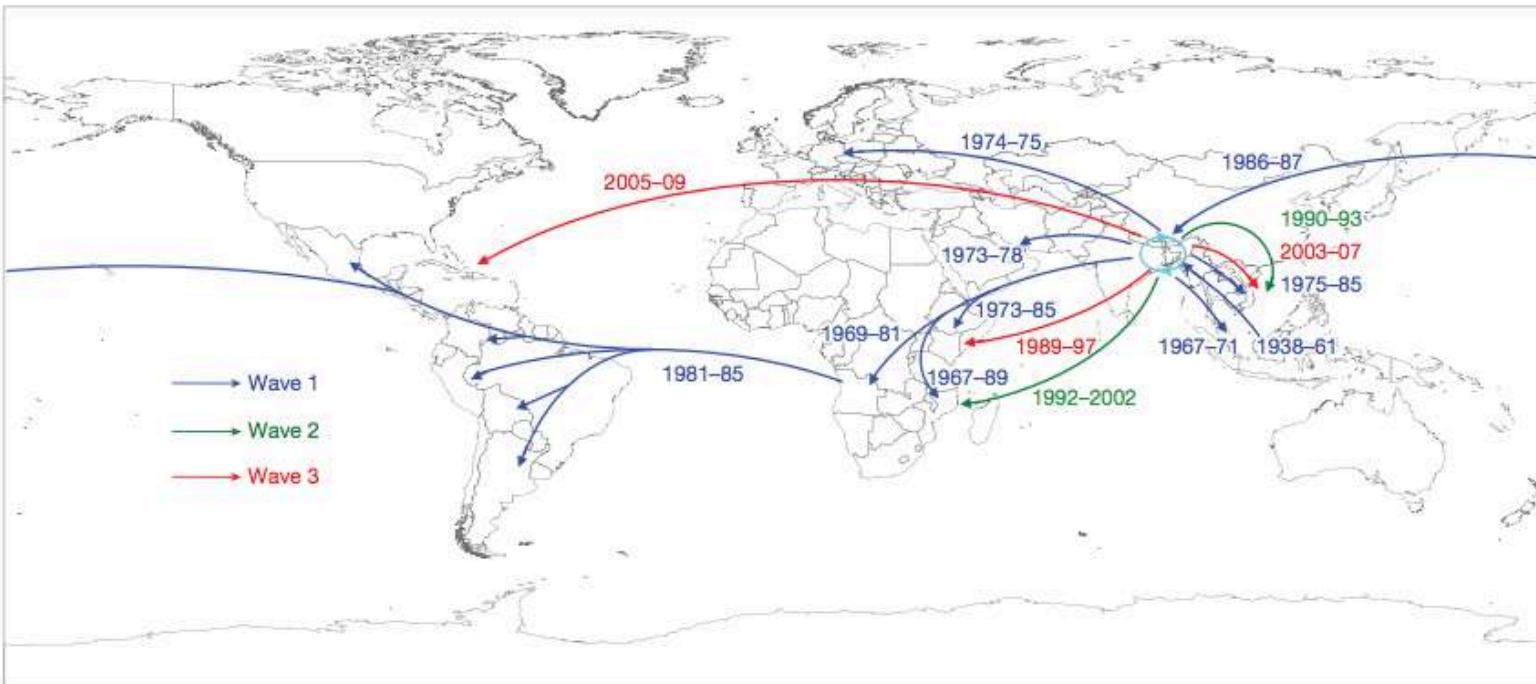


Novembre et al Nature (2008)



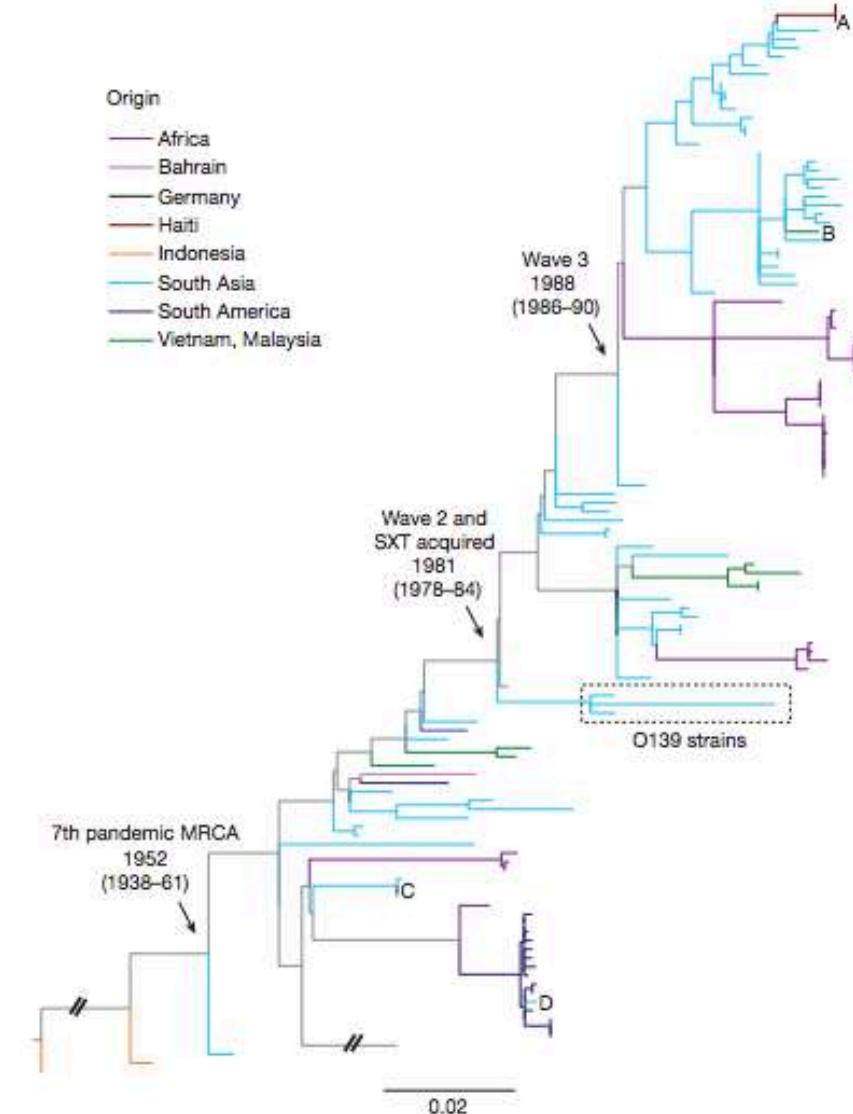
http://www.genomenext.com/casestudies_post/population-scale-analysis-genomic-samples-analyzed-from-2504-individuals-in-1-week/

Population genomics



Origin

- Africa
- Bahrain
- Germany
- Haiti
- Indonesia
- South Asia
- South America
- Vietnam, Malaysia



2013

8000家庭破碎 聯合國遭控傳染霍亂

2013-10-11 by : 阿嫲

5725 ❤️ f t

一直以來，聯合國給世人的印象多是促進世界永續發展的正面印象，但對海地居民來說，聯合國卻成了當地人最恐懼的劊子手，最新報導就指出，因聯合國駐軍而散佈的霍亂已經造成8,000人死亡。

BBC綜合報導，聯合國派駐的維和部隊(UN peacekeepers)意外將細菌帶到海地境內，在當地造成霍亂大流行，自2010年爆發至今，霍亂已經在海地造成8,000人病死，這也讓海地成為目前全世界霍亂疫病最嚴重的地區。

聯合國是兇手

儘管許多調查指出聯合國就是霍亂源頭，但海地數度請願要求補償未果，現在海地的代表律師團就上訴紐約法院，控告聯合國是造成海地霍亂疫情的元凶。

2016

聯合國坦承：我們將霍亂帶進了海地

2016-08-19 by : 泥仔

15040 ❤️ f t

將近六年的時間，聯合國終於承認海地的霍亂疫情與他們有關。到目前為止，已經有數十萬名居民感染上霍亂、一萬名海地人因霍亂而去世。

維和部隊惹的禍？

由於海地過去都沒有類似霍亂症狀的疾病，部分專家也發現海地的霍亂細菌種類與尼泊爾的種類是一樣的，因此懷疑是聯合國在尼泊爾的維和部隊將霍亂弧菌帶進海地。但將近六年來，聯合國一直都否認這樣的指控。

聯合國坦承與疫情爆發有關

在本周三(17)，聯合國副發言人哈奇(Farhan Haq)聲明：「過去幾年來，聯合國有鑑於海地初期的瘟疫爆發與我們有些關係，聯合國決定要多做些什麼。」他也強調聯合國會在接下來兩個月內有所行動。

UK launches whole genome sequence alliance to map spread of coronavirus

The Wellcome Sanger Institute will collaborate with expert groups across the country to analyse the genetic code of COVID-19 samples circulating in the UK, providing public health agencies with a unique tool to combat the virus

COVID-19 Genomics UK Consortium - comprised of the NHS, Public Health Agencies, Wellcome Sanger Institute, and numerous academic institutions - will deliver large scale, rapid sequencing of the cause of the disease and share intelligence with hospitals, regional NHS centres and the Government.

Samples from patients with confirmed cases of COVID-19 will be sent to a network of sequencing centres which currently includes Belfast, Birmingham, Cambridge, Cardiff, Edinburgh, Exeter, Glasgow, Liverpool, London, Norwich, Nottingham, Oxford and Sheffield.

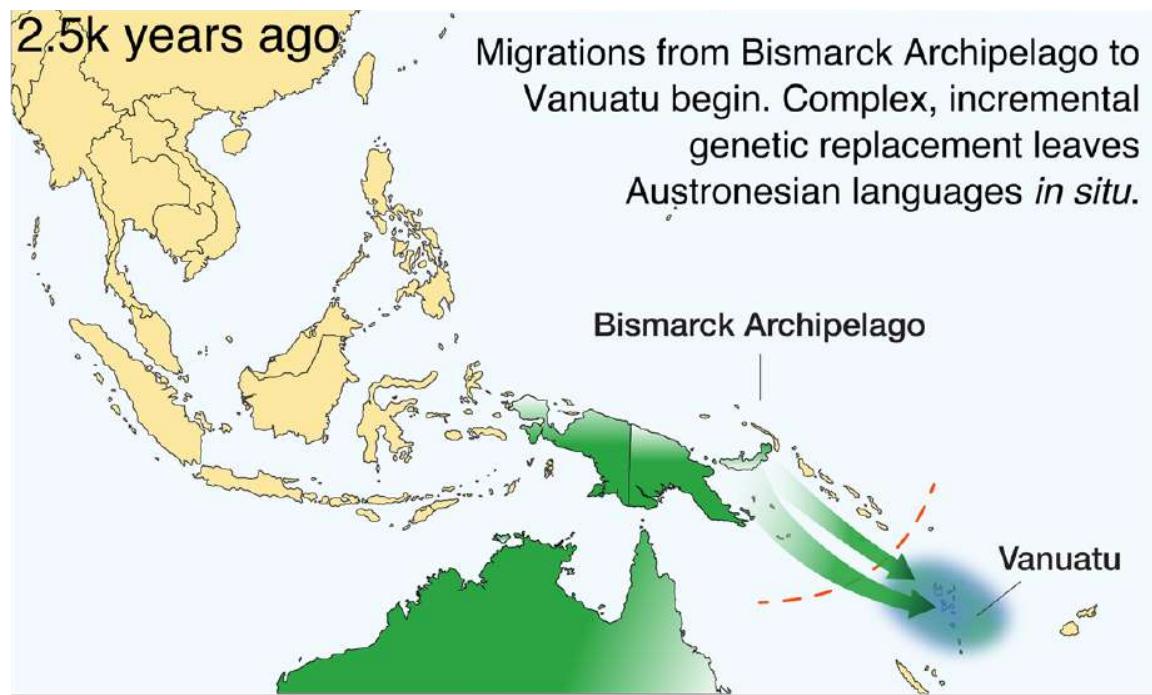
The Wellcome Sanger Institute, one of the world's most advanced centres of genomes and data, will collaborate with expert groups across the country to analyse the genetic code of COVID-19 samples circulating in the UK and in doing so, give public health agencies and clinicians a unique, cutting-edge tool to combat the virus.

By looking at the whole virus genome in people who have had confirmed cases of COVID-19, scientists can monitor changes in the virus at a national scale to understand how the virus is spreading and whether different strains are emerging. This will help clinical care of patients and save lives.

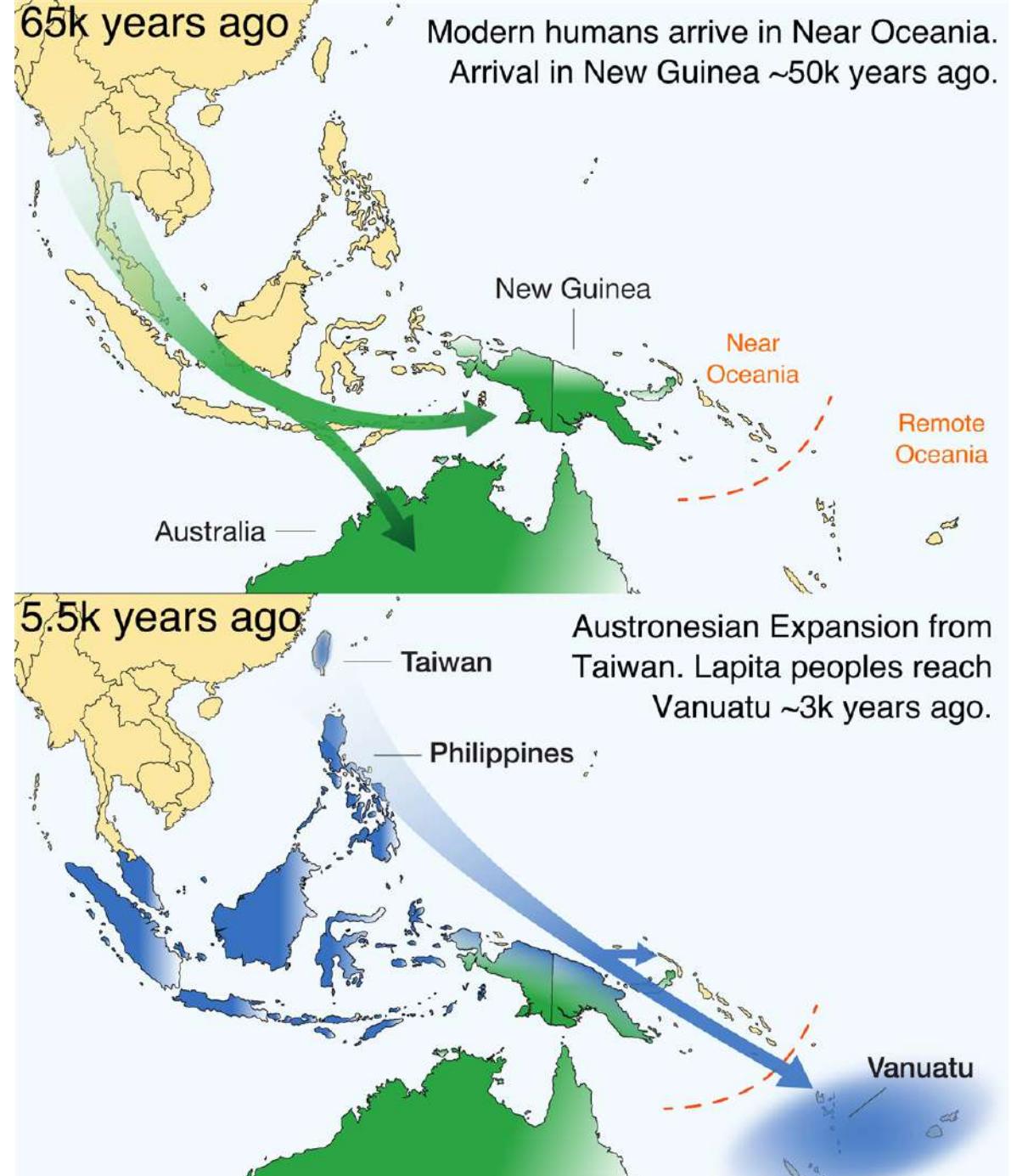
<https://www.ncbi.nlm.nih.gov/pubmed/28102248>

Population genomics

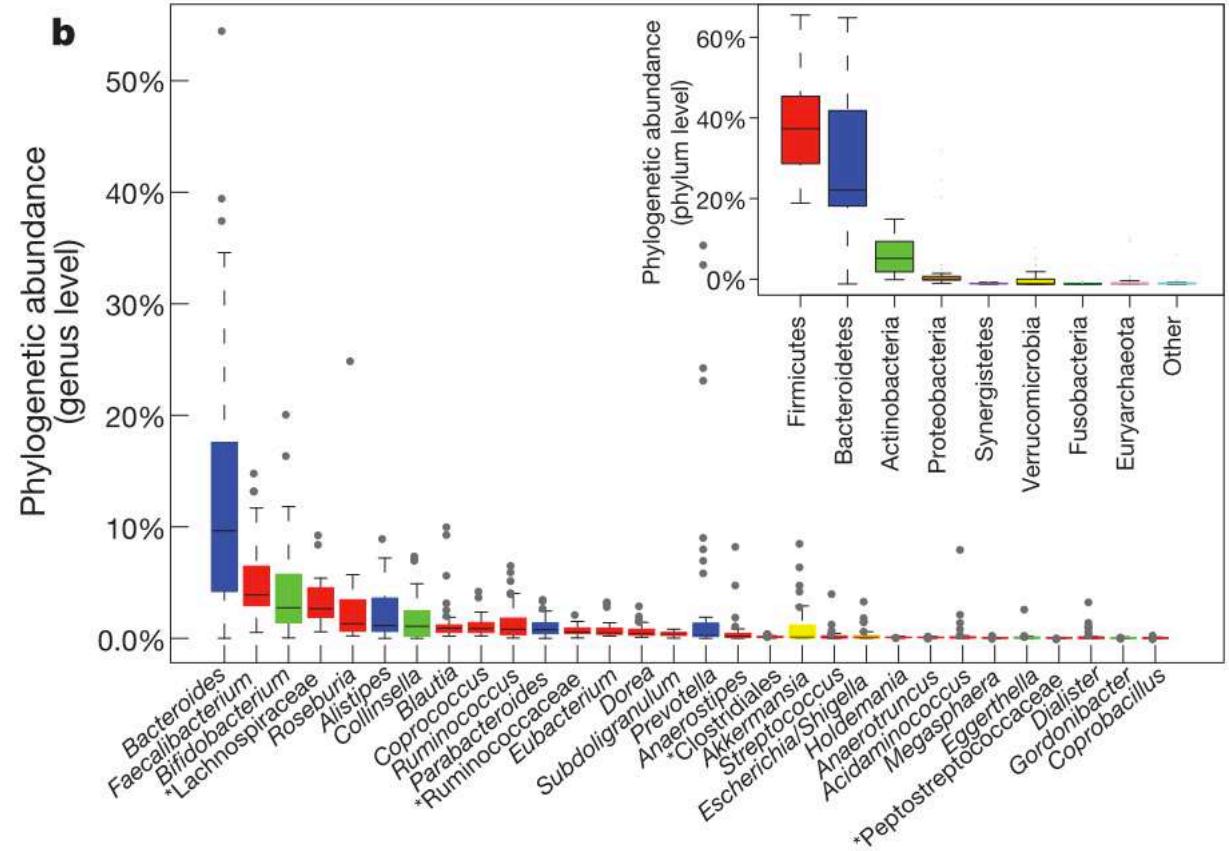
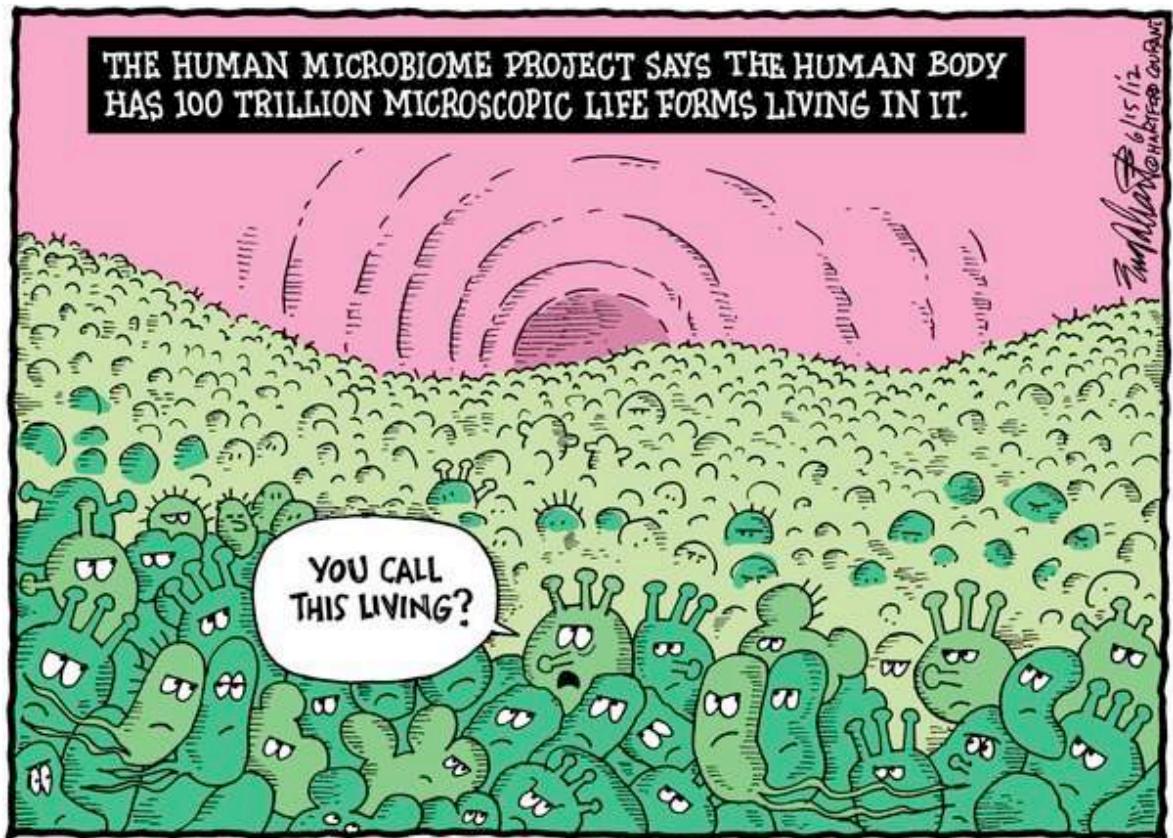
There is a puzzling mismatch in the population history of the Pacific: why do most people across Remote Oceania (the vast area stretching from Vanuatu to Easter Island-Rapa Nui) speak languages of the Austronesian language family that expanded into this region only 3,000 years ago, yet carry a component of genetic ancestry from a much older source population in Near Oceania (the area including New Guinea and its surrounding islands, see the top of Figure 1)?



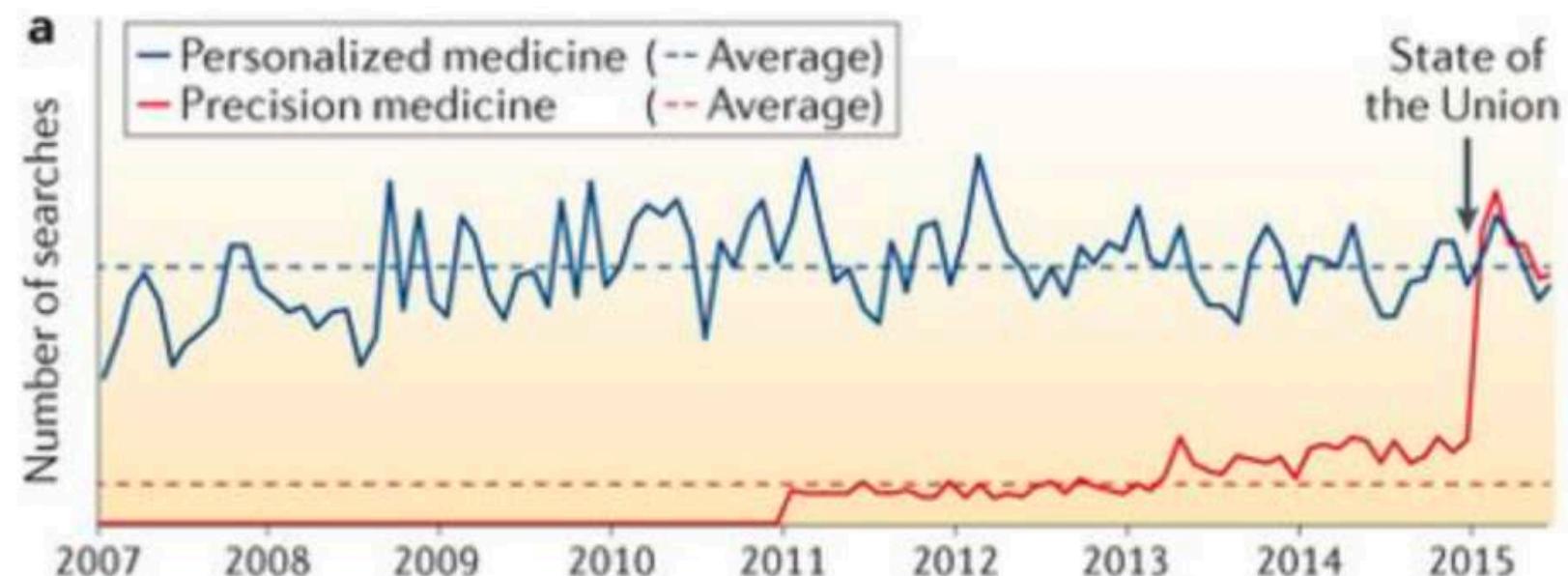
Credit: Hans Sell, MPI-SHH, adapted from Skoglund et al. 2016 Nature



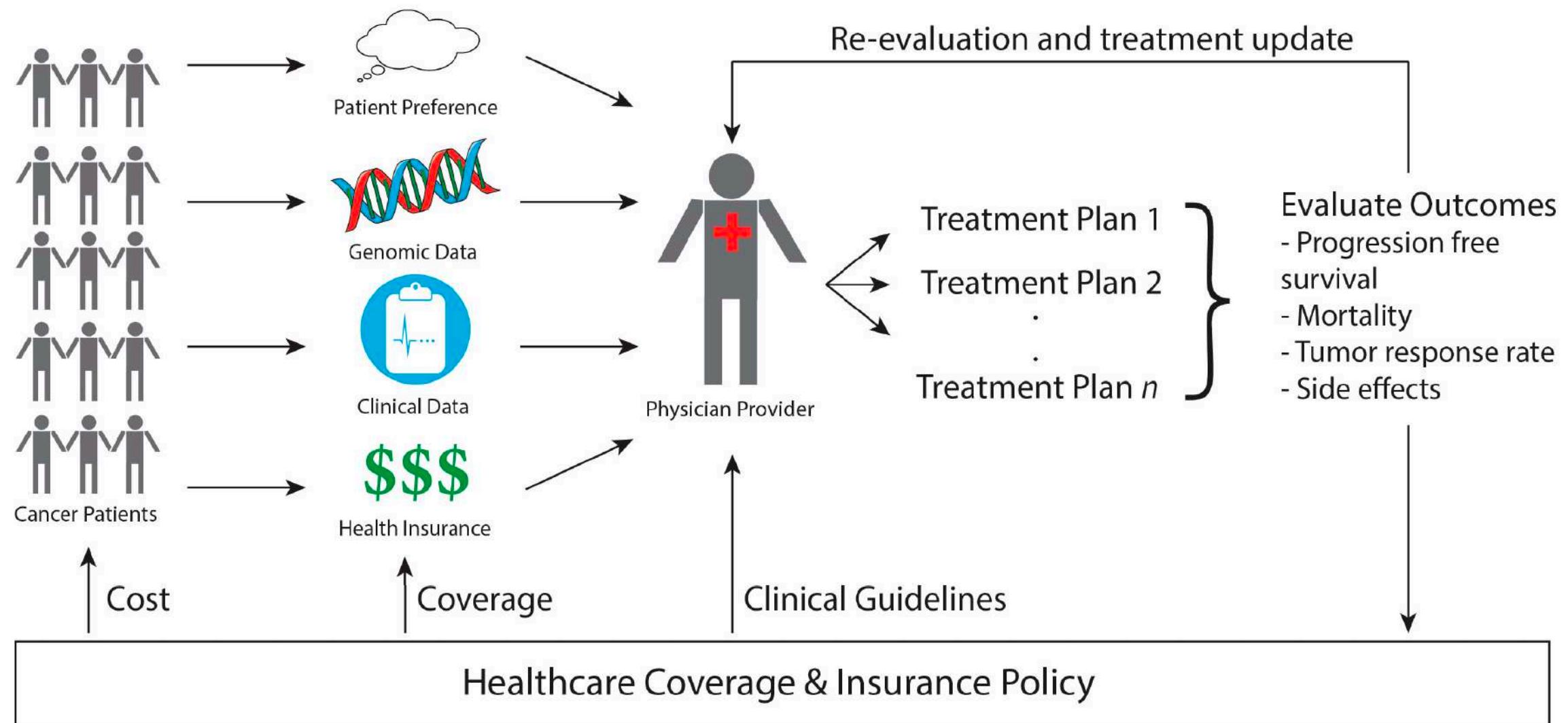
Metagenomics



Precision medicine



Outline of precision medicine



PM examples

Table 1 | Examples of precision medicine

Condition	Gene	Action
Mendelian disease		
Cystic fibrosis	<i>CFTR</i>	Specific therapies such as ivacaftor and a combination of lumacaftor and ivacaftor
Long QT syndrome	<i>KCNQ1, KCNH2 and SCN5A</i>	Specific therapy for patients with <i>SCN5A</i> mutations
Duchenne muscular dystrophy	<i>DMD</i>	Ongoing phase III clinical trials of exon-skipping therapies
Malignant hyperthermia susceptibility	<i>RYR1</i>	Avoid volatile anaesthetic agents; avoid extremes of heat
Familial hypercholesterolaemia (FH)	<i>PCSK9, APOB and LDLR</i>	<ul style="list-style-type: none"> Heterozygous FH (HeFH): eligible for PCSK9 inhibitor drugs Homozygous FH (HoFH): eligible for PCSK9 inhibitor drugs in addition to lomitapide and mipomersen
Dopa-responsive dystonia	<i>SPR</i>	Therapy with dopamine precursor L-dopa and the serotonin precursor 5-hydroxytryptophan
Thoracic aortic aneurysm	<i>SMAD3, ACTA2, TGFBR1, TGFBR2 and FBN1</i>	Customization of surgical thresholds based on patient genotype
Left ventricular hypertrophy	<i>MYH7, MYBPC3, GLA and TTR</i>	Sarcomeric cardiomyopathy, Fabry disease and transthyretin cardiac amyloid disease have specific therapies
Precision oncology		
Lung adenocarcinoma	<i>EGFR and ALK</i>	Targeted kinase inhibitors, such as gefitinib and crizotinib
Breast cancer	<i>HER2</i>	HER2 (also known as ERBB2)-targeted treatment, such as trastuzumab and pertuzumab
Gastrointestinal stromal tumour	<i>KIT</i>	Targeted KIT kinase activity inhibitors, such as imatinib
Melanoma	<i>BRAF</i>	BRAF inhibitors, such as vemurafenib and dabrafenib
Pharmacogenomics		
Warfarin sensitivity	<i>CYP2C9 and VKORC1</i>	Adjust dosage of warfarin or consider alternative anticoagulant
Clopidogrel sensitivity, post-stent procedure	<i>CYP2C19</i>	Consider alternative antiplatelet therapy (for example, prasugrel or ticagrelor)
Thiopurine sensitivity	<i>TPMT</i>	Reduce thiopurine dosage or consider alternative agent
Codeine sensitivity	<i>CYP2D6</i>	Avoid use of codeine; consider alternatives such as morphine and non-opioid analgesics
Simvastatin sensitivity	<i>SLCO1B1</i>	Reduce dose of simvastatin or consider an alternative statin; consider routine creatine kinase surveillance

Summary of outcomes in Oncology PM Studies

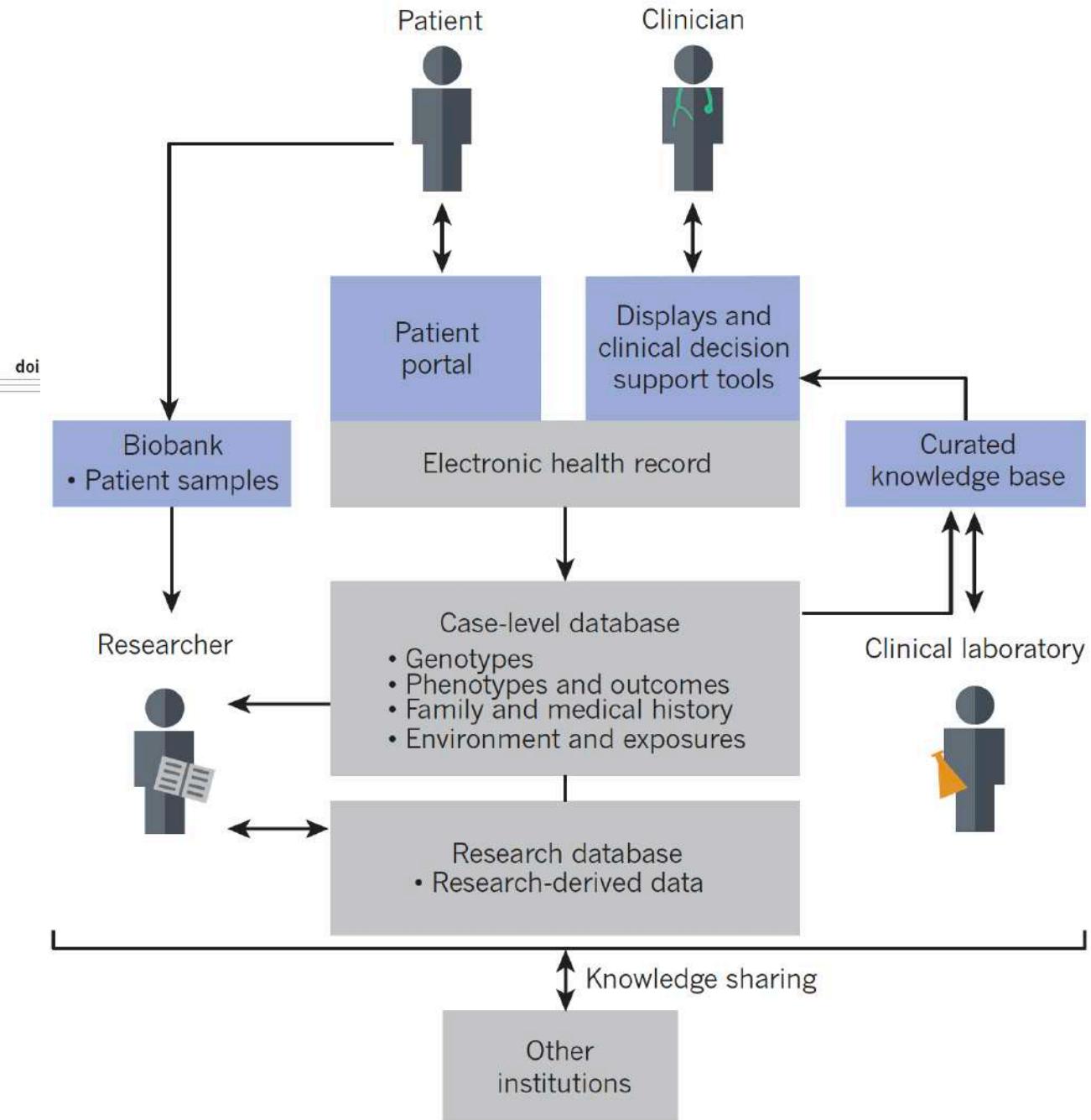
Study	Sample Size	Most Prevalent Tumor Types	Outcomes Reported
Tsimberidou et al. <i>Clin. Cancer Res.</i> 2012 [5]	291 patients with one molecular aberration (175 treated with matched therapy, 116 control)	Colorectal, melanoma, lung, ovarian	Matched group had improved ORR (27% vs. 5%), TTF (median 5.2 vs. 2.2 month), OS (median 13.4 vs. 9.0 month)
Radovich et al. <i>Oncotarget</i> 2016 [6]	101 patients with sequencing and follow up (44 treated with matched therapy, 57 control)	Soft tissue sarcoma, breast, colorectal	Matched group had improved PFS (86 vs. 49 days)
Schwaederle et al. <i>Mol. Cancer Ther.</i> 2016 [7]	180 patients with sequencing and follow up (87 treated with matched therapy, 93 control)	Gastrointestinal, breast, brain	Matched group had improved PFS (4.0 vs. 3.0 month), TRR (34.5% vs. 16.1% achieving SD/PR/CR)
Kris et al. <i>JAMA</i> 2014 [8]	578 patients with oncogenic driver and followup (260 with matched therapy, 318 control)	Lung only	Matched group had improved survival (median 3.5 vs. 2.4 years)
Aisner et al. <i>J. Clin. Oncol.</i> 2016 [9]	187 patients with targetable alteration and follow up (112 with matched therapy, 74 control)	Lung only	Matched group had improved survival (median 2.8 vs. 1.5 years)
Stockley et al. <i>Genome Med.</i> 2016 [10]	245 patients with sequencing matched to clinical trials (84 on matched trial, 161 control)	Gynecological, lung, breast	Matched group had improved ORR (19% vs. 9%)
LeTourneau et al. <i>Lancet Oncol.</i> 2015 [11]	RCT with 195 patients with molecular aberration (99 treated with matched therapy, 96 control)	Gastrointestinal, breast, brain	No difference in PFS between groups

ORR = overall response rate, TTF = time to treatment failure, OS = overall survival, PFS = progression free survival, TRR = tumor response rate, SD = stable disease, PR = partial response, CR = complete response, RCT = randomized controlled trial. Matched group indicates patients matched to a therapy based on sequencing results.

REVIEW

Building the foundation for genomics in precision medicine

Samuel J. Aronson^{1,2} & Heidi L. Rehm^{1,3,4,5}



Building the foundation for genomics in precision medicine

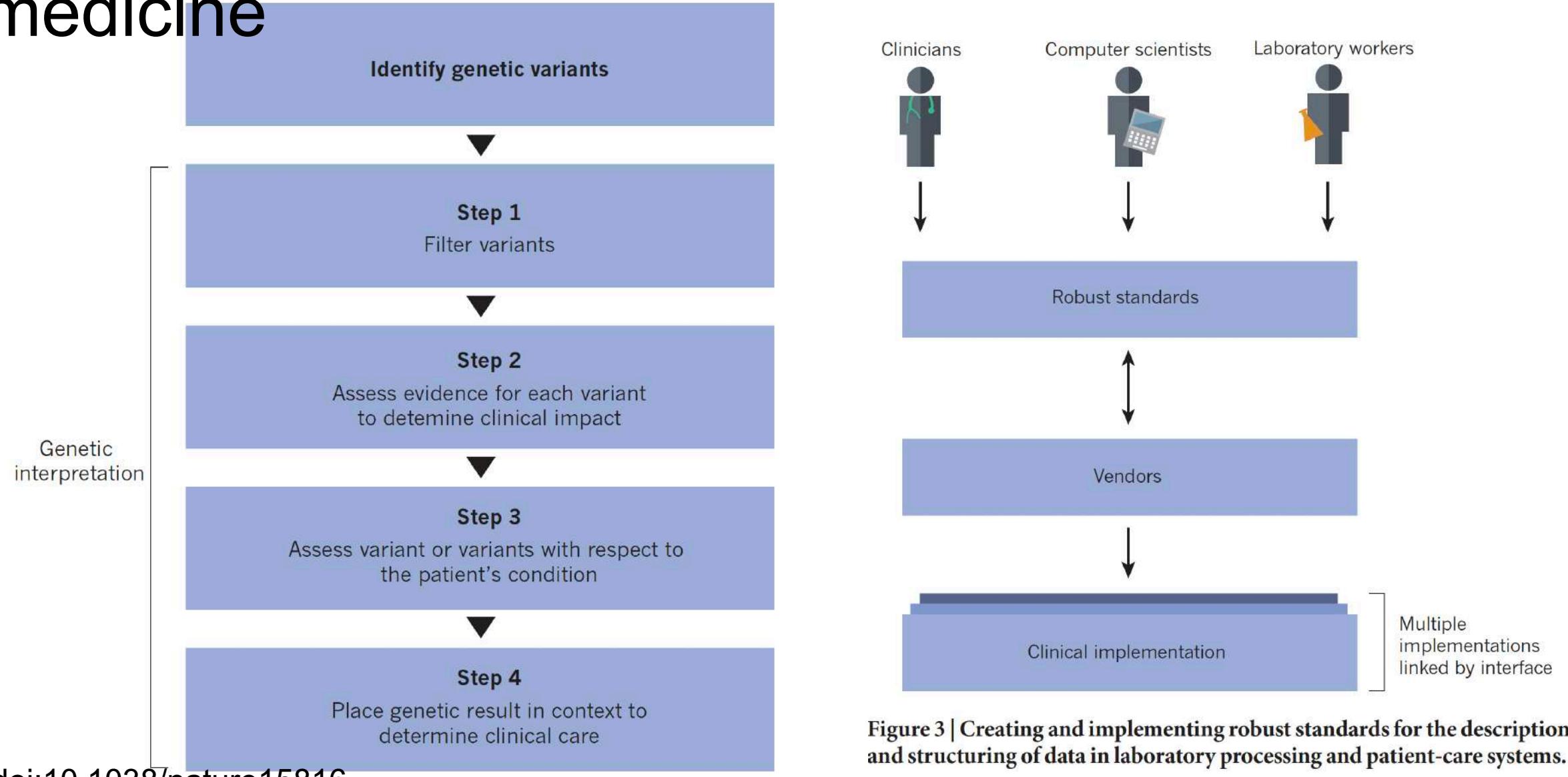
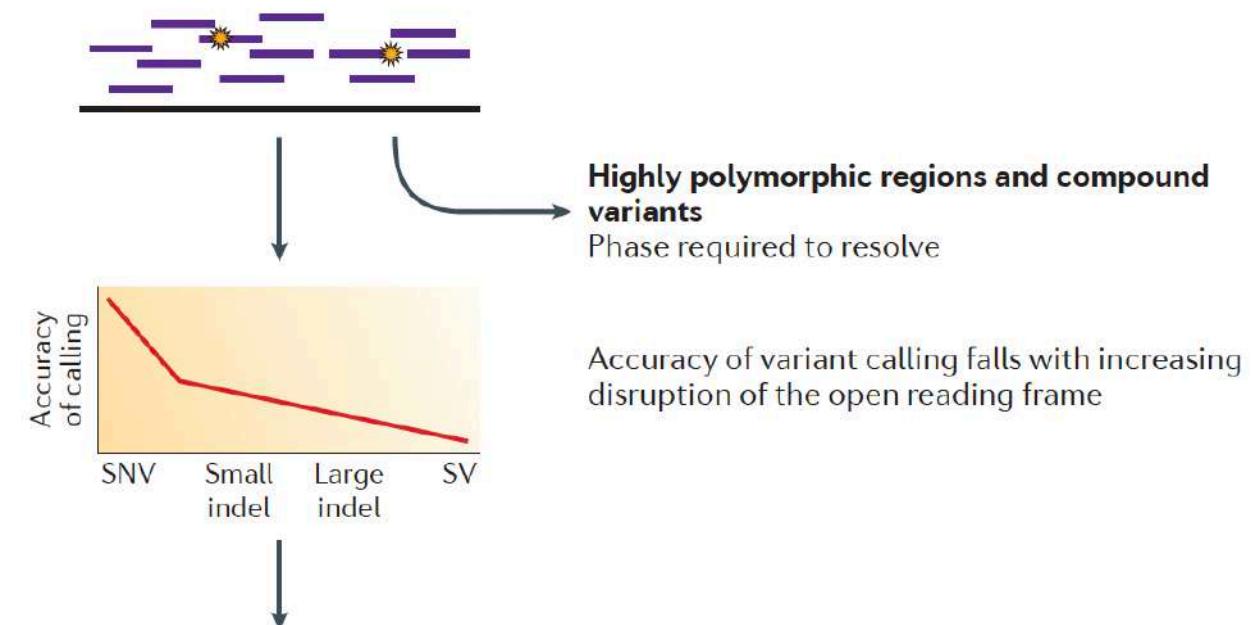
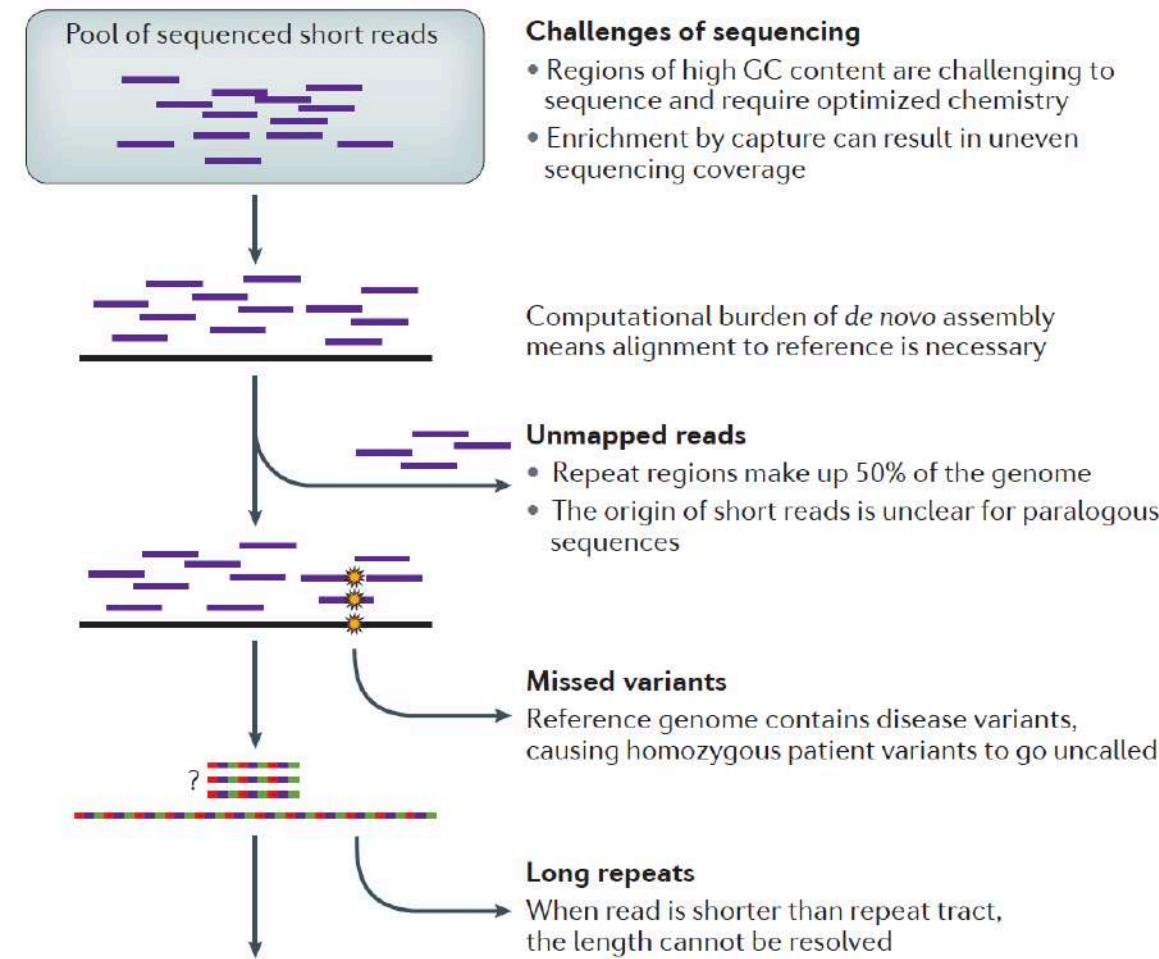
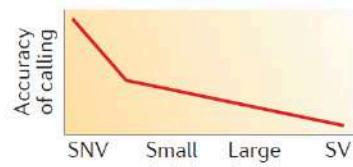


Figure 3 | Creating and implementing robust standards for the description and structuring of data in laboratory processing and patient-care systems.

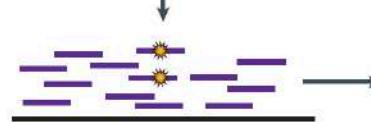
Challenges and reduced accuracies



Challenges and reduced accuracies



Accuracy of variant calling falls with increasing disruption of the open reading frame



Position	REF	ALT	Call
Chr14:23,456,332	T	A	0/1

Final VCF file

- File of appropriately called variants
- The VCF should contain a call at every position or patients homozygous for risk alleles present in the reference will be missed

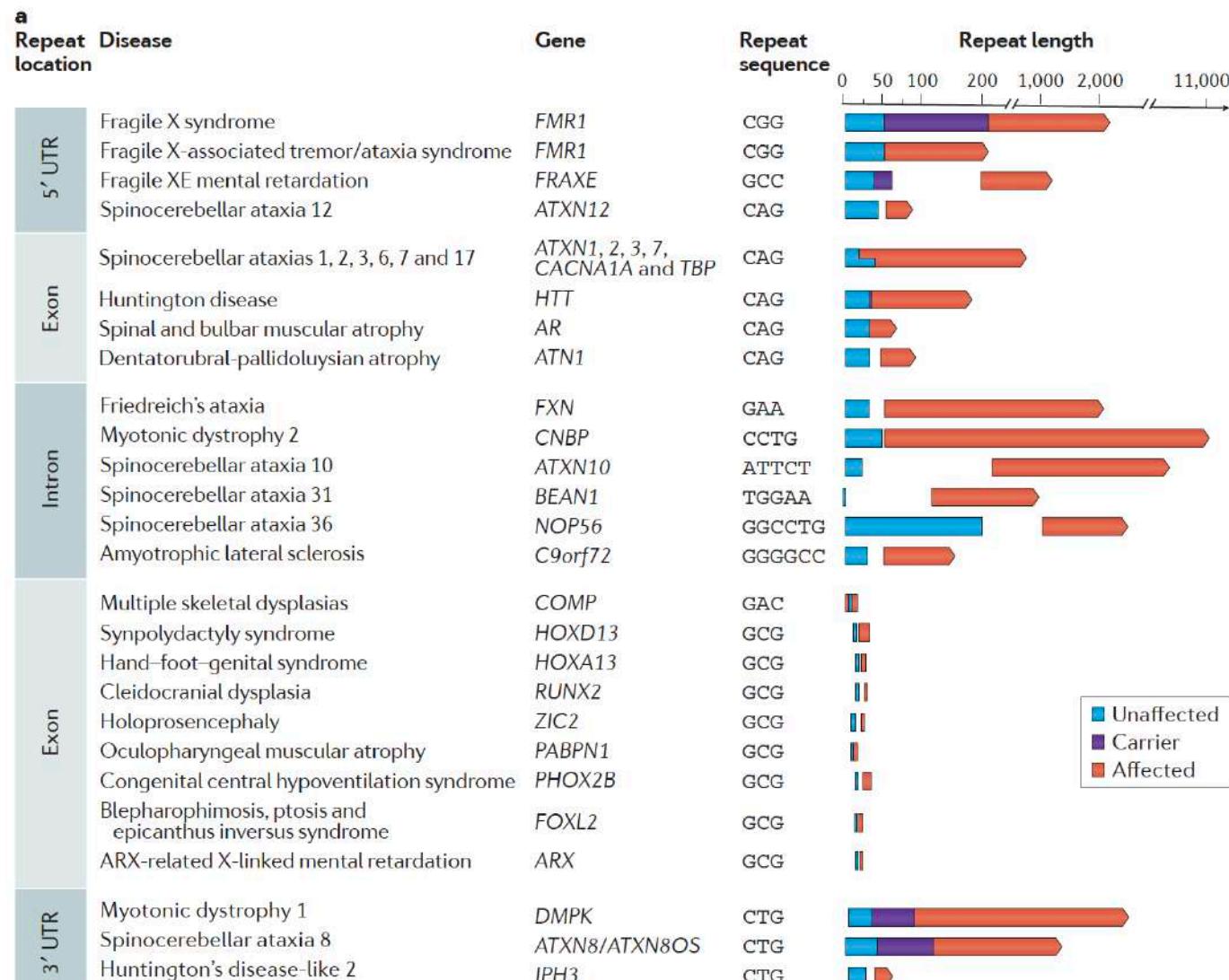
Position	REF	ALT	Call
Chr14:23,456,332	T	A	0/1
Chr14:23,678,972	C	G	1/1
...			

Variants filtered based on standard metrics, such as population frequency and known disease-associated genes

Position	REF	ALT	Call
Chr14:23,456,332	T	A	0/1
Chr14:23,678,972	C	G	1/1
...			

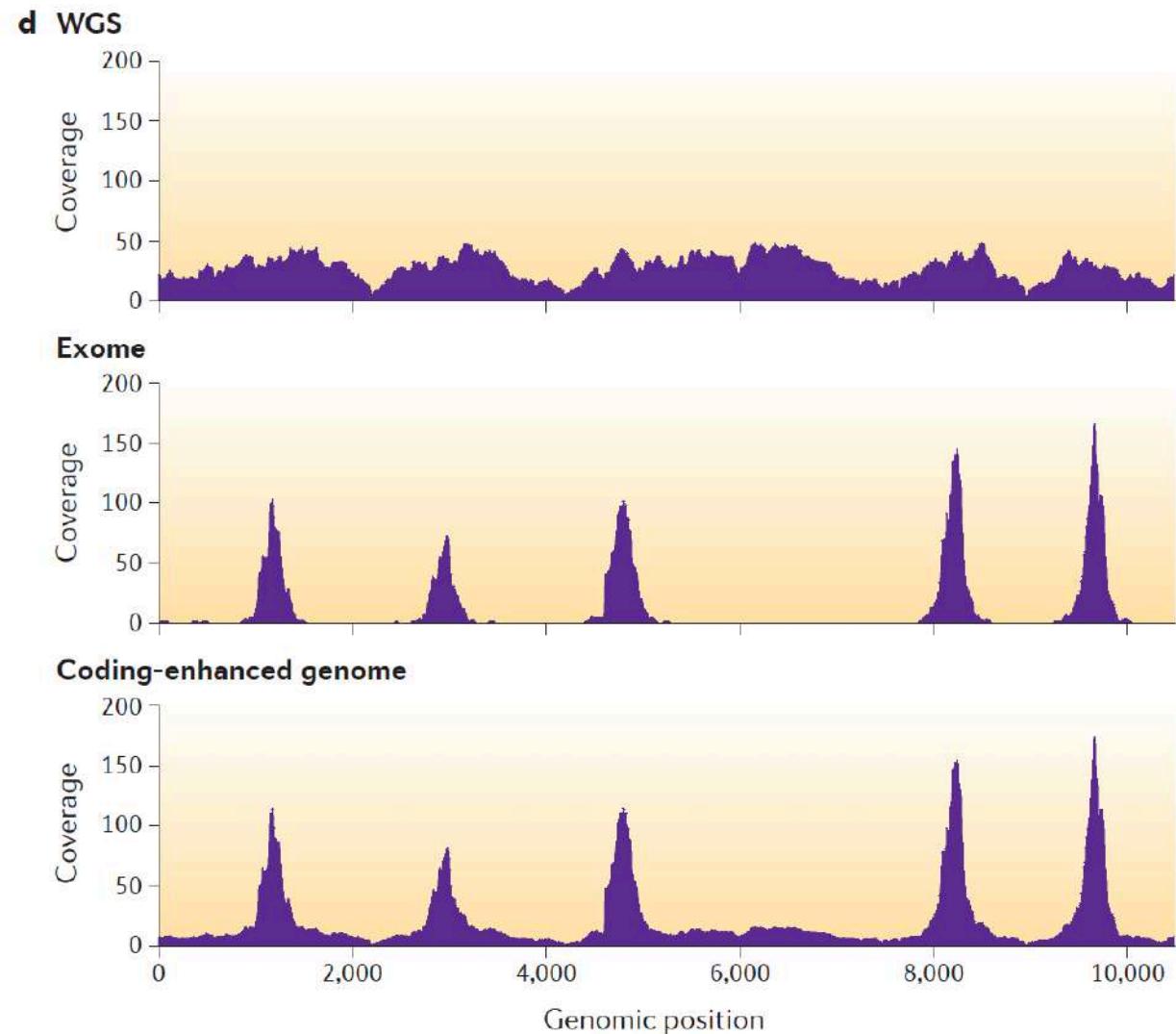
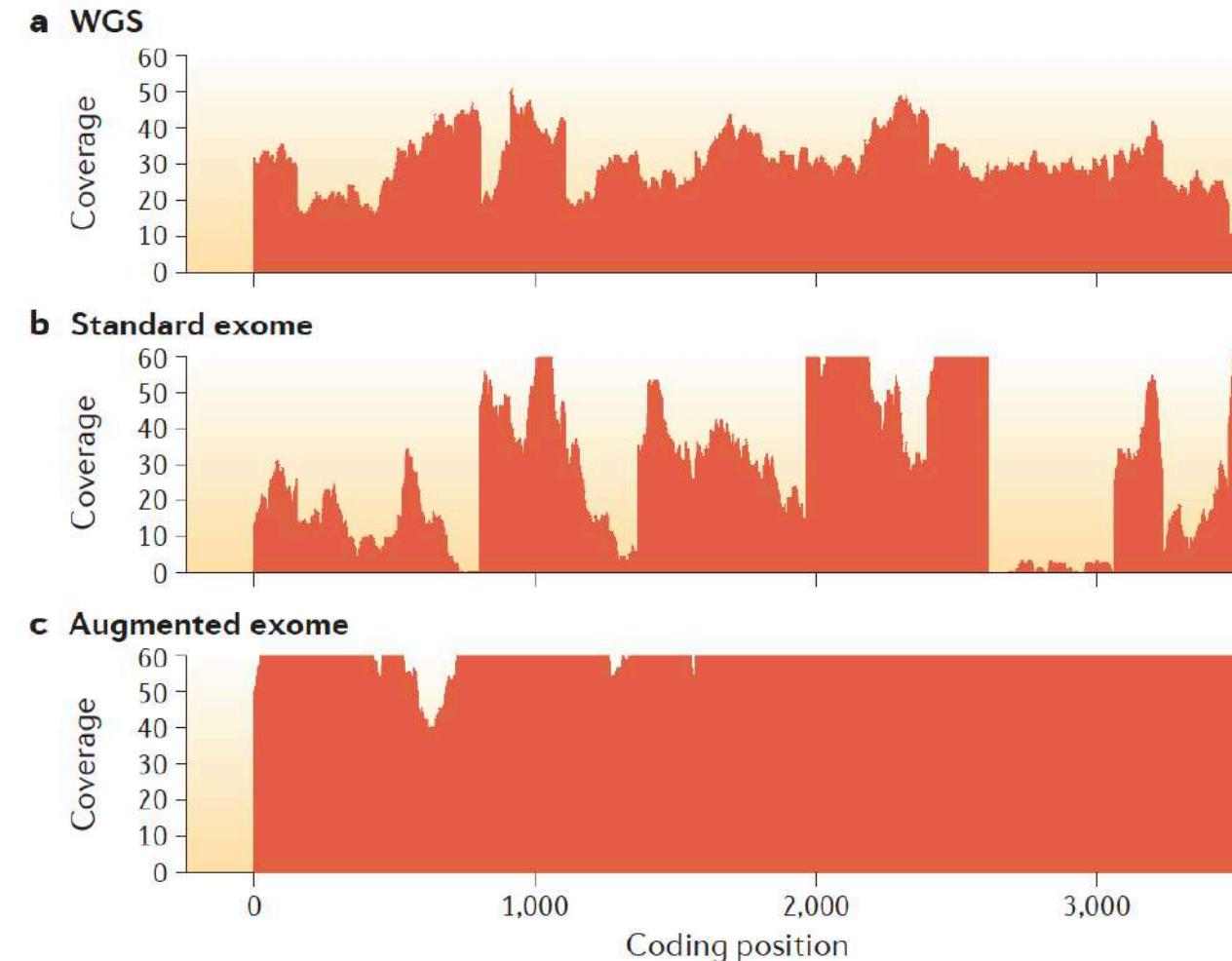
Causality determined by magnitude and dependency of effect

Downstream treatment and disease management are influenced by knowledge of disease-causing gene and variant



Resolution

Zoomed out



Transcriptomics / RNAseq

Applications of RNAseq

Discovery / Annotation

- Find new genes
- Find new transcripts
- Find new ncRNAs, xxx, xxx
- Gene fusion

Comparison / Quantification : given X conditions, find the effect of Y on

- **expression**
- **Isoform abundance, splice patterns, transcript boundaries**

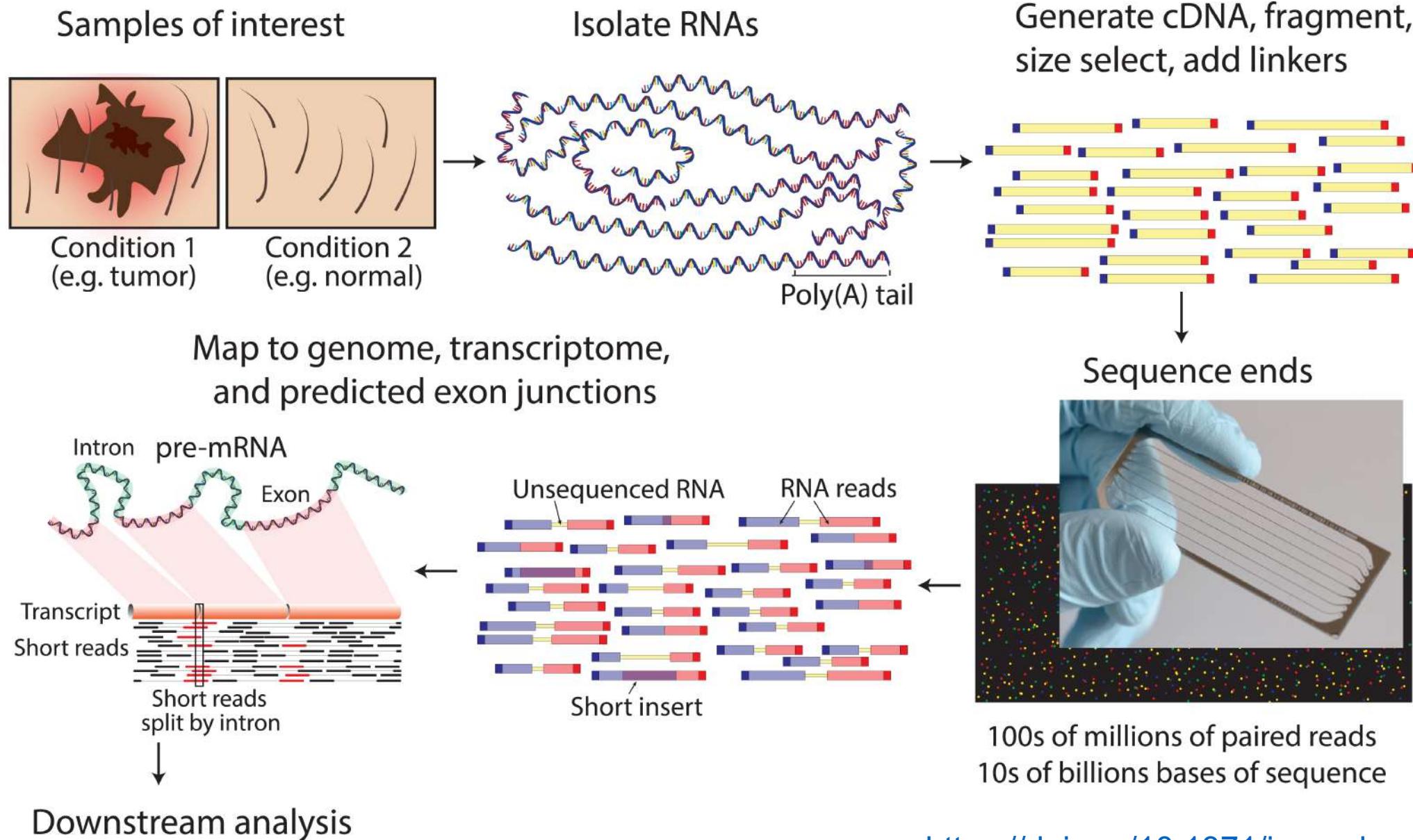
Introduction to differential gene expression analysis using RNA-seq

Written by Friederike Dündar, Luce Skrabaneck, Paul Zumbo

September 2015
updated March 20, 2018

<http://chagall.med.cornell.edu/RNASEQcourse/Intro2RNaseq.pdf>

RNA-seq data generation

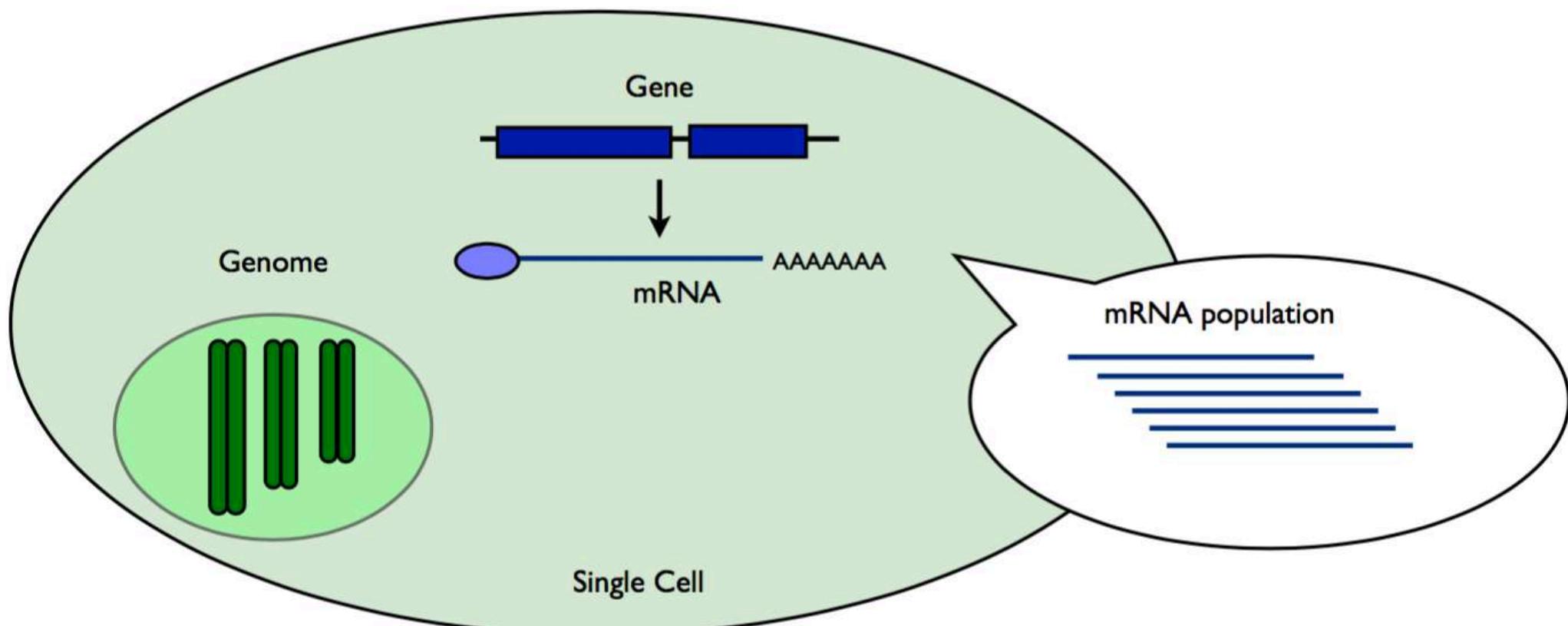


Types of experiments

Transcriptome Complexity:

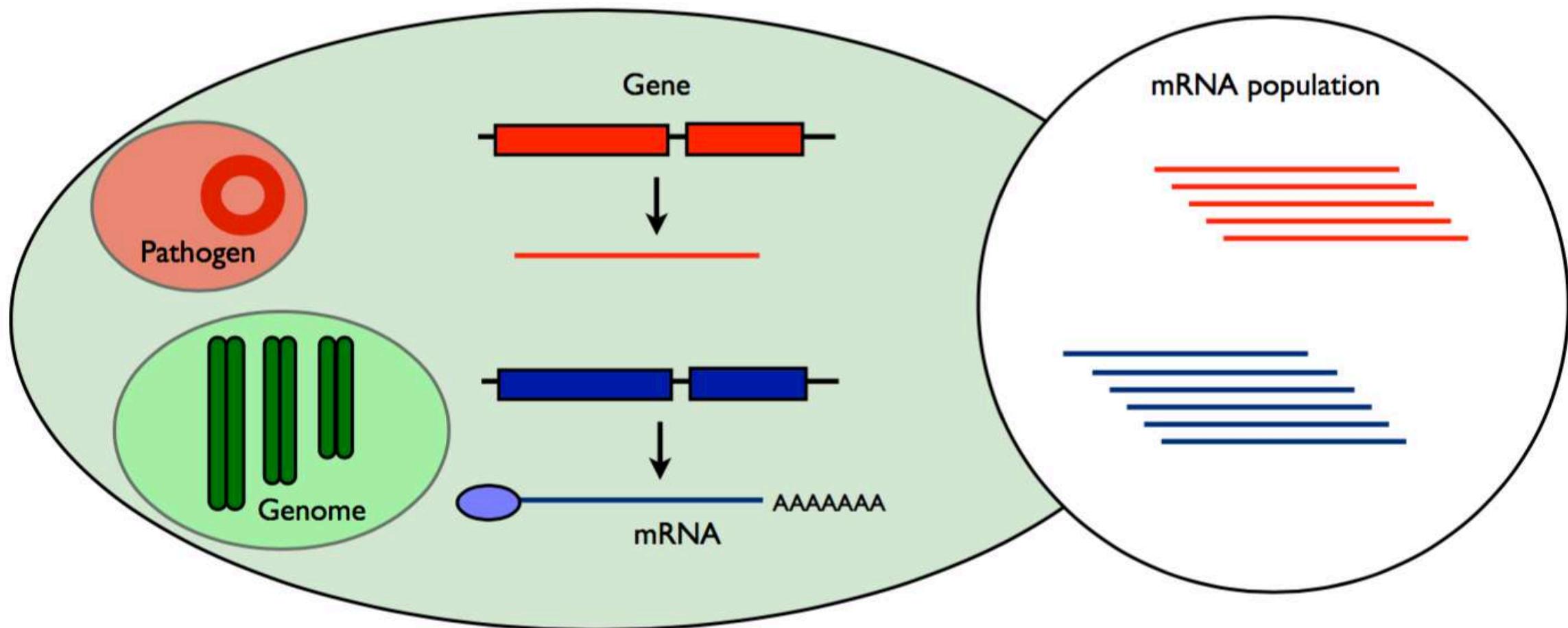
Simple System:

One Genome => Gene 1 copy => Single mRNA



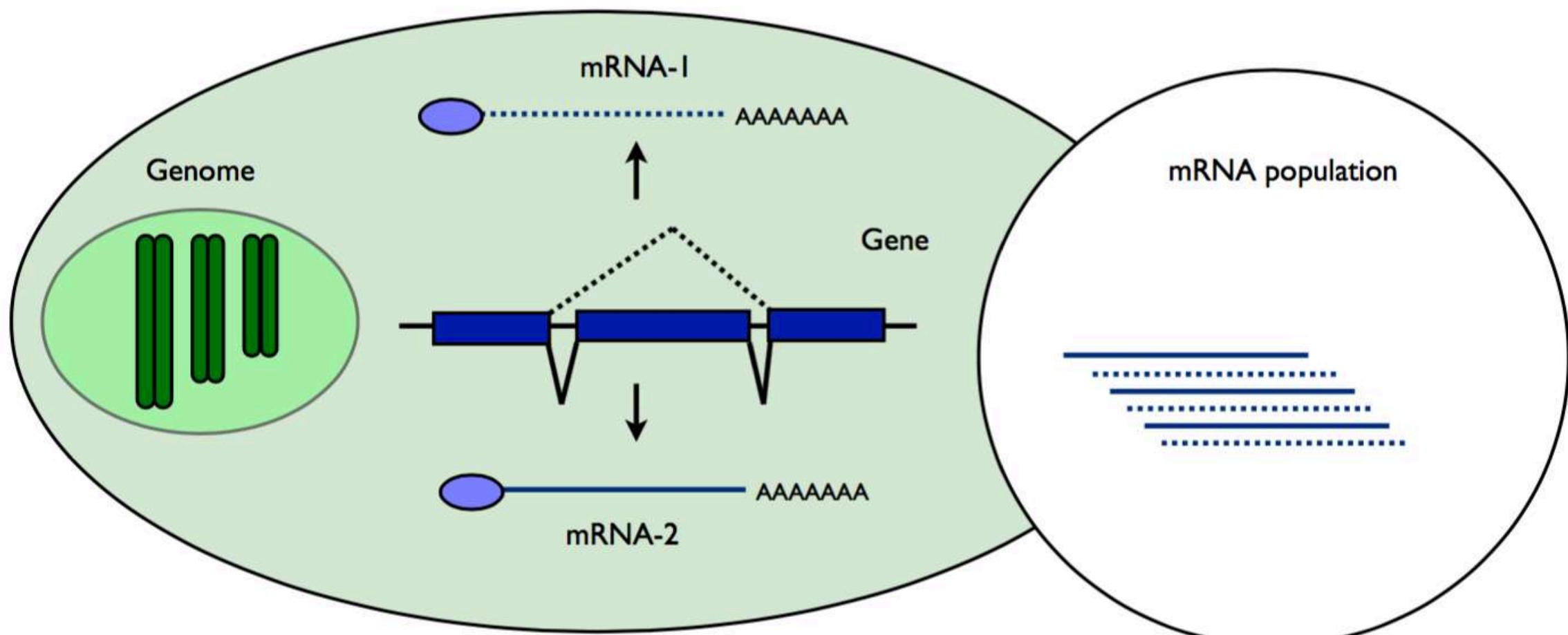
How many species we are analyzing ?

- 1) Problems to isolate a single species (rhizosphere)
- 2) Species interaction study (plant-pathogen)



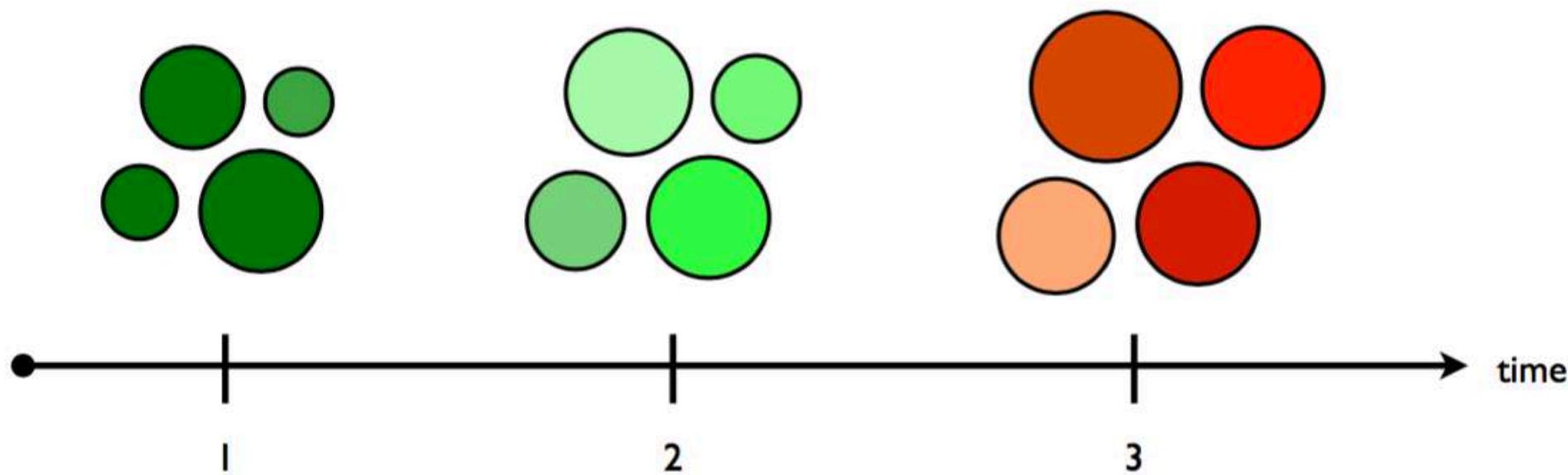
How many isoforms we expect for each allele ?

1) Alternative splicings



Is the study performed at different time points?

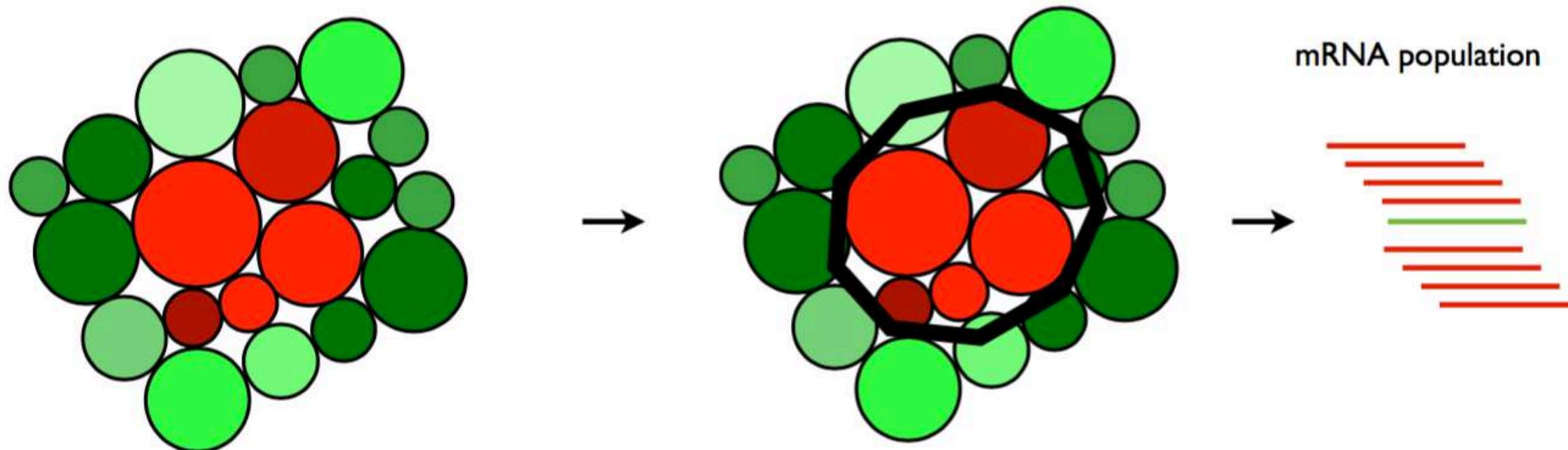
- 1) Developmental stages (difficult to select the same)**
- 2) Response to a treatment**



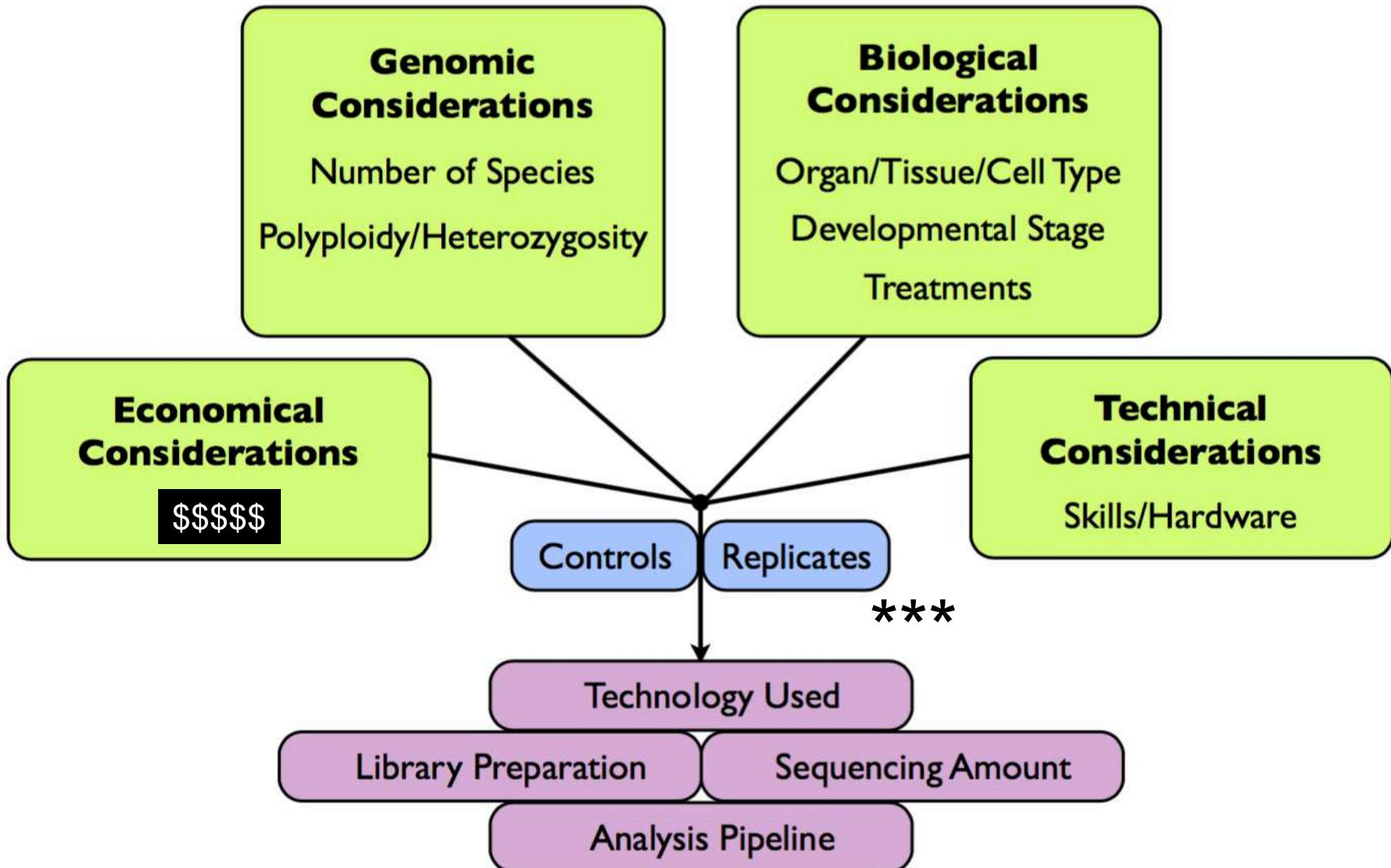
Is the study performed with different parts?

- 1) Organ specific**
- 2) Tissue/Cell type specific**

(Laser Capture Microdissection, LCM)



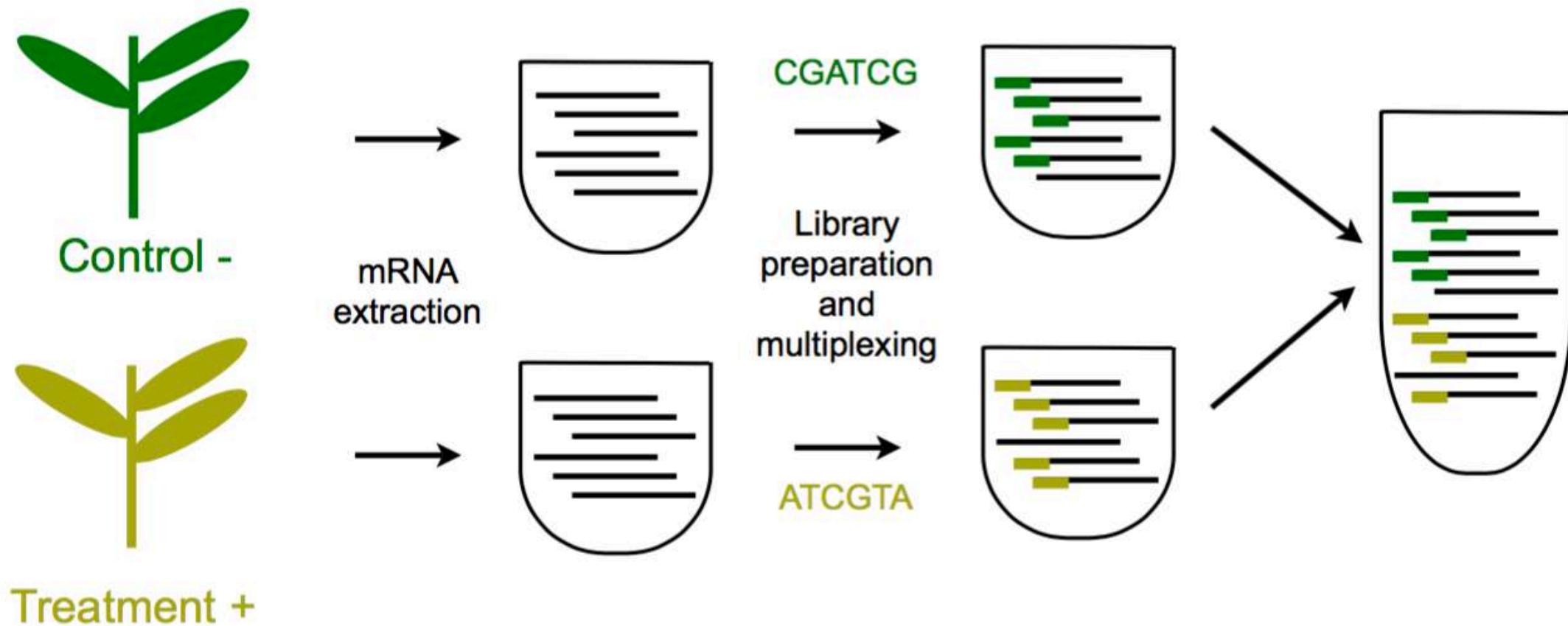
Experimental design



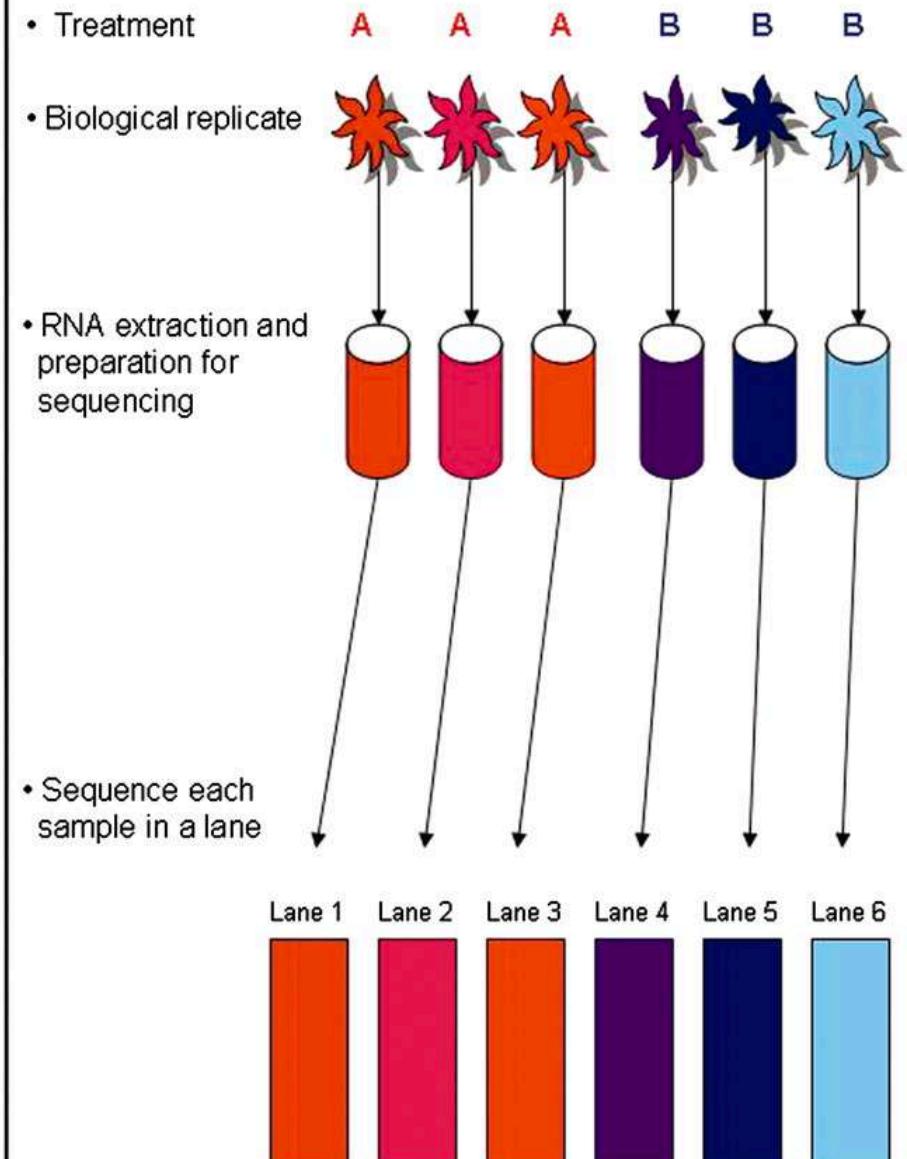
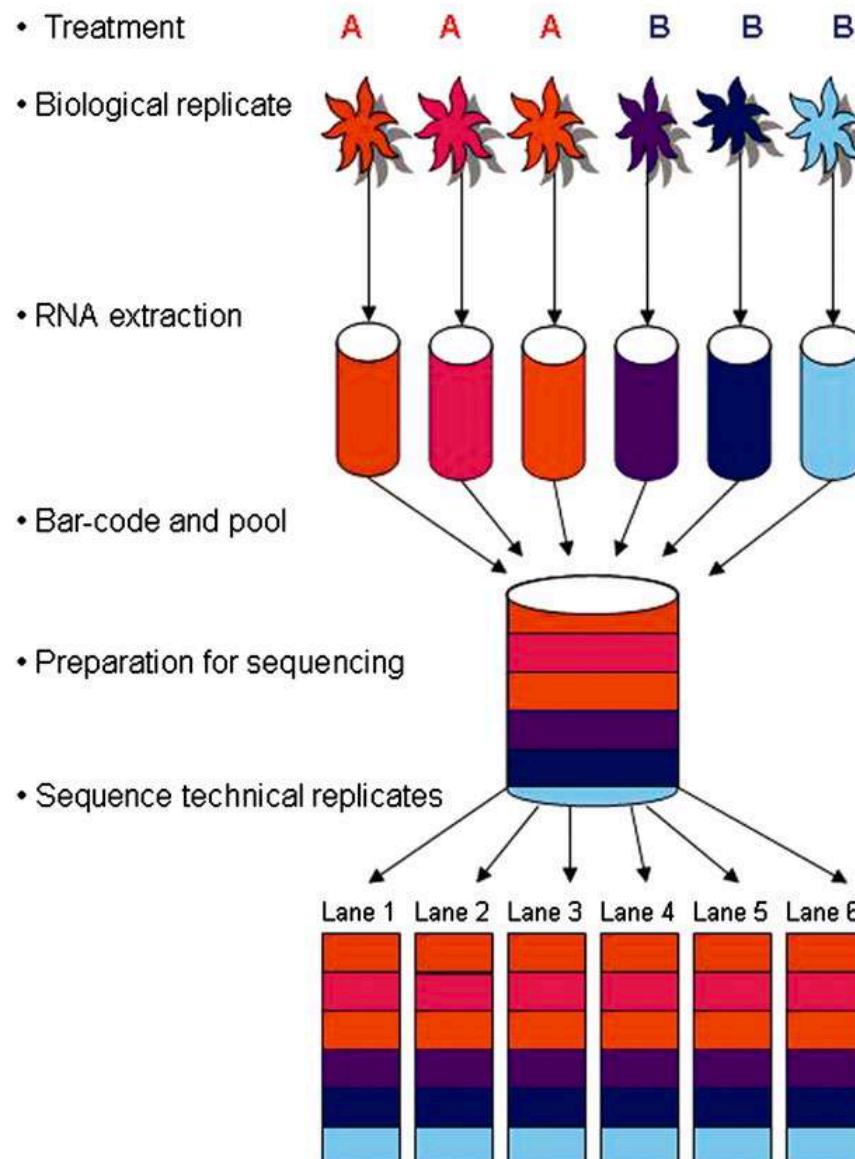
Prep and treatment

Sequencing of multiple samples can be performed using **multiplexing**.

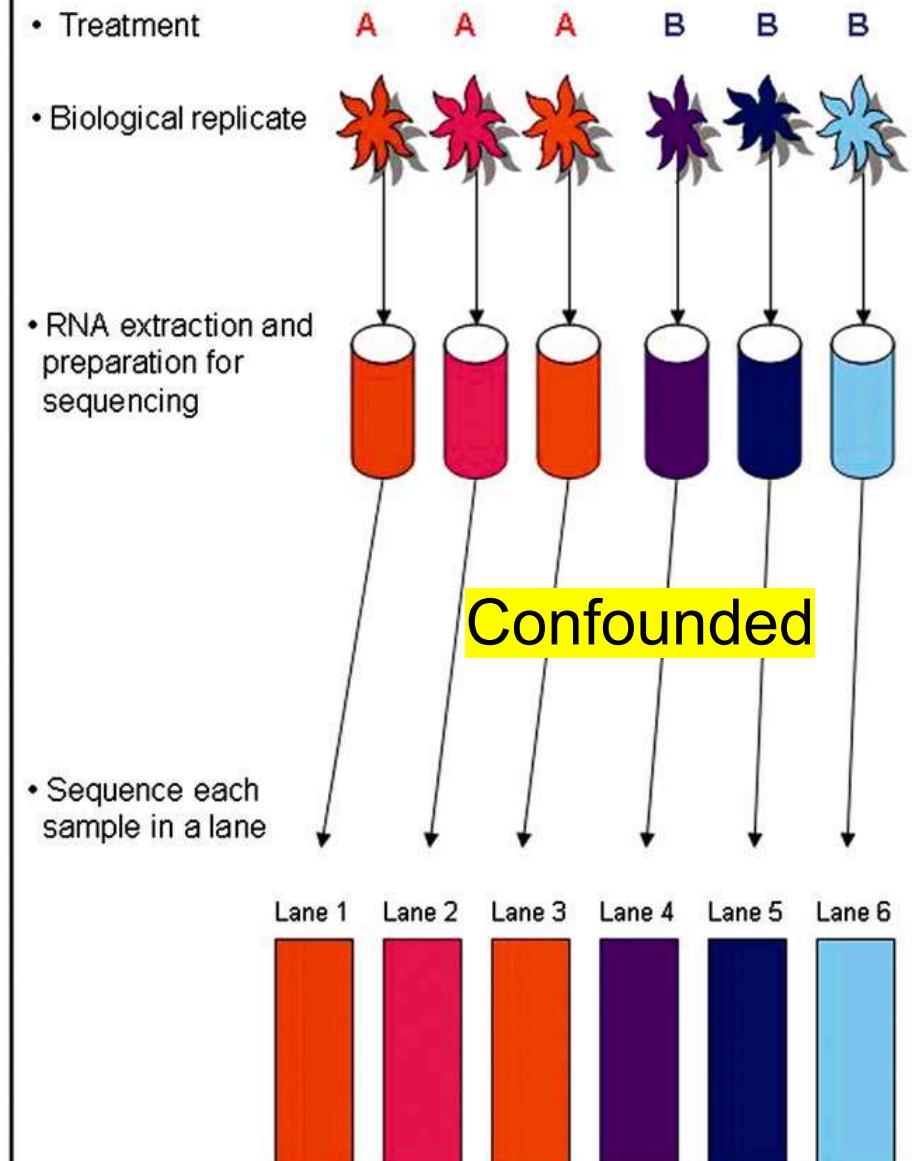
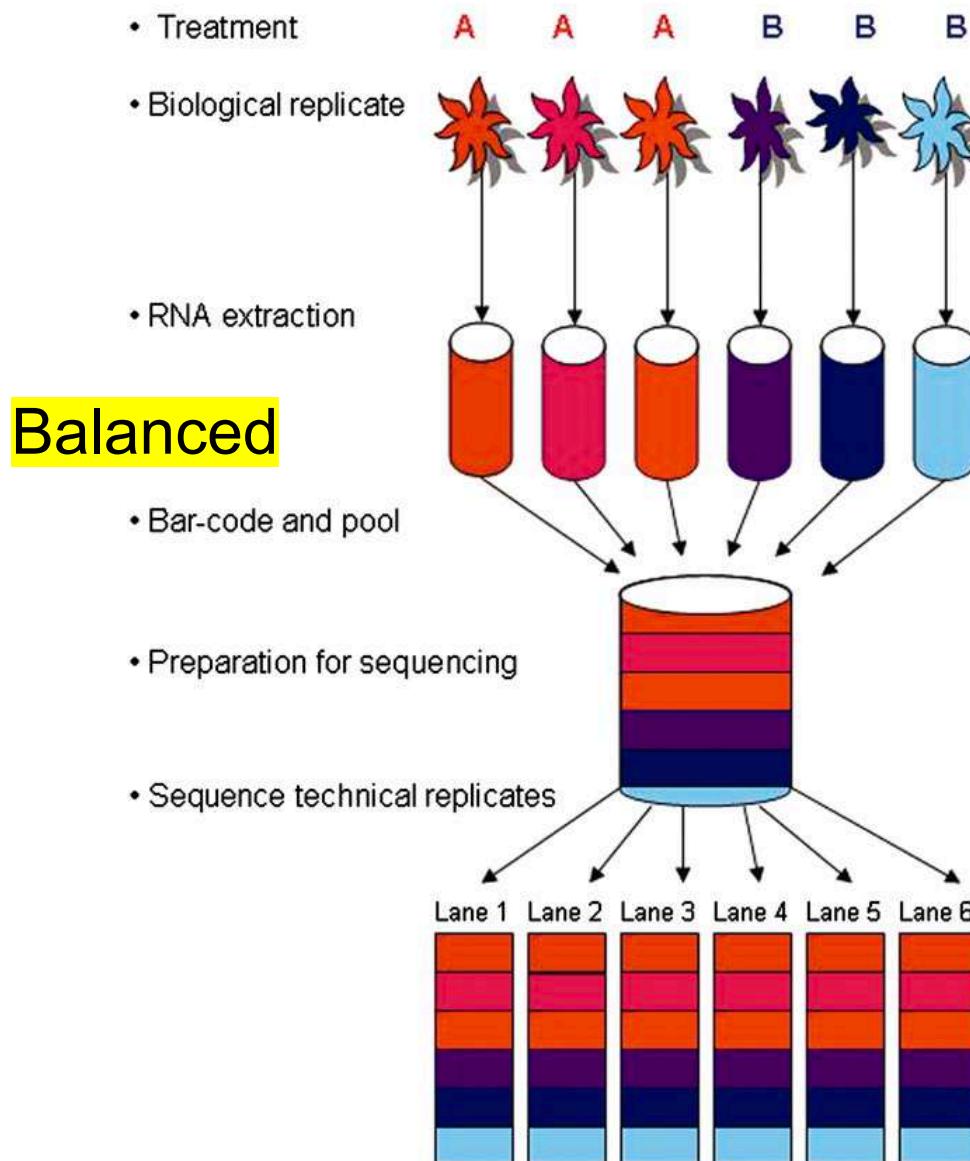
The multiplexing add a tag/**barcode** of 4-6 nucleotides during the library preparation to identify the sample. Common kits can add up to 96 different tags.



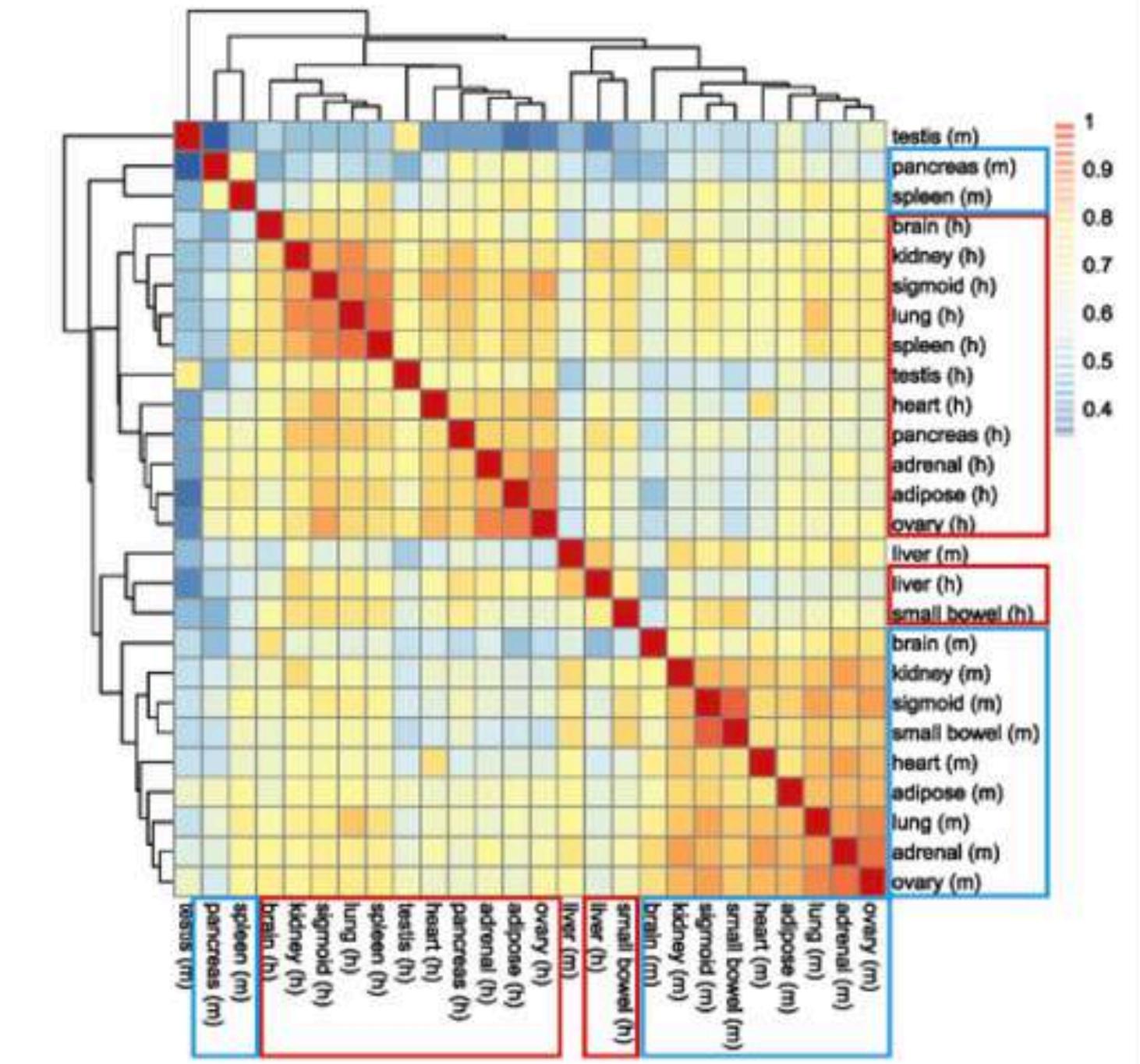
Which of the following designs is correct?



Which of the following designs is correct?



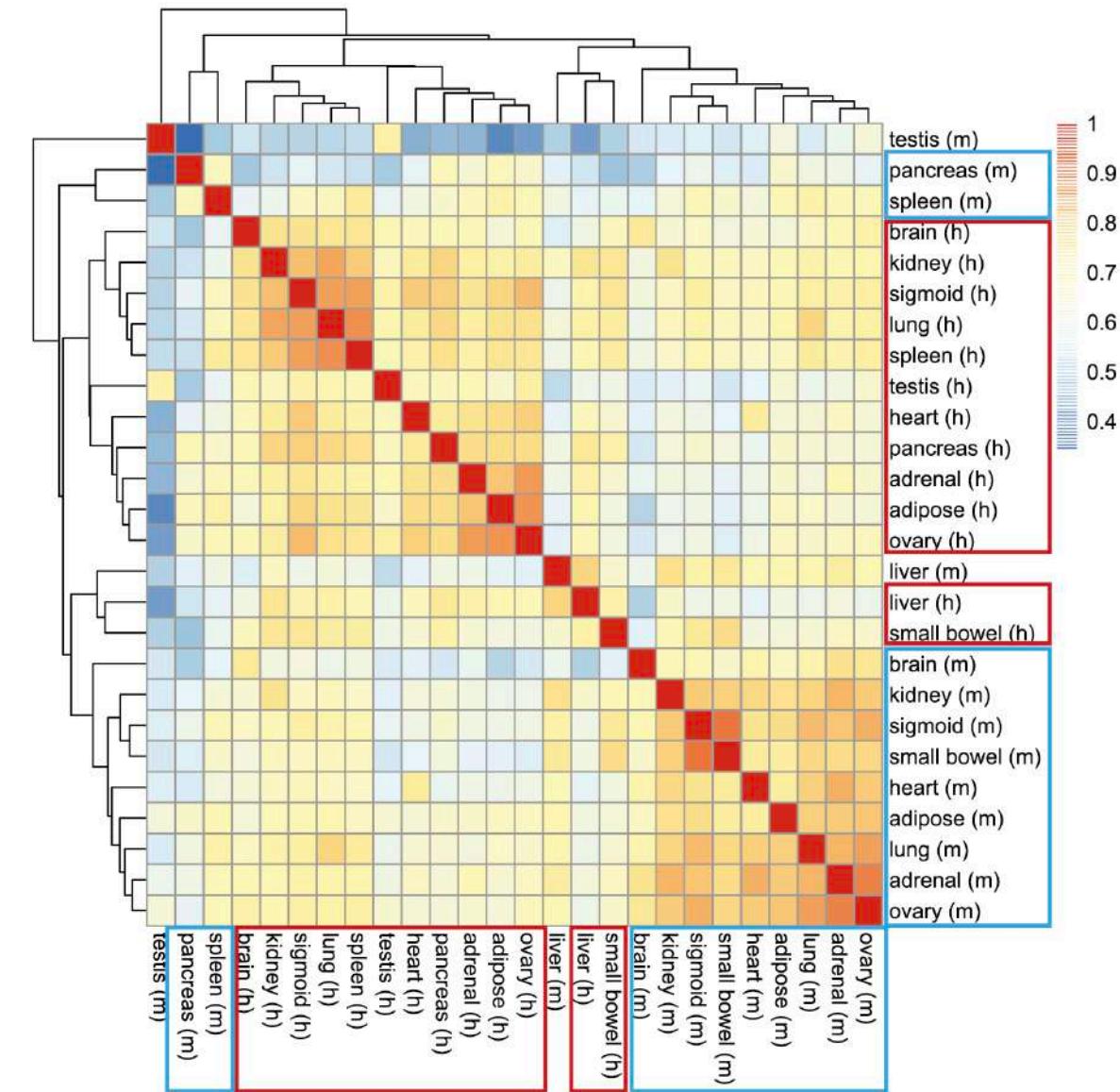
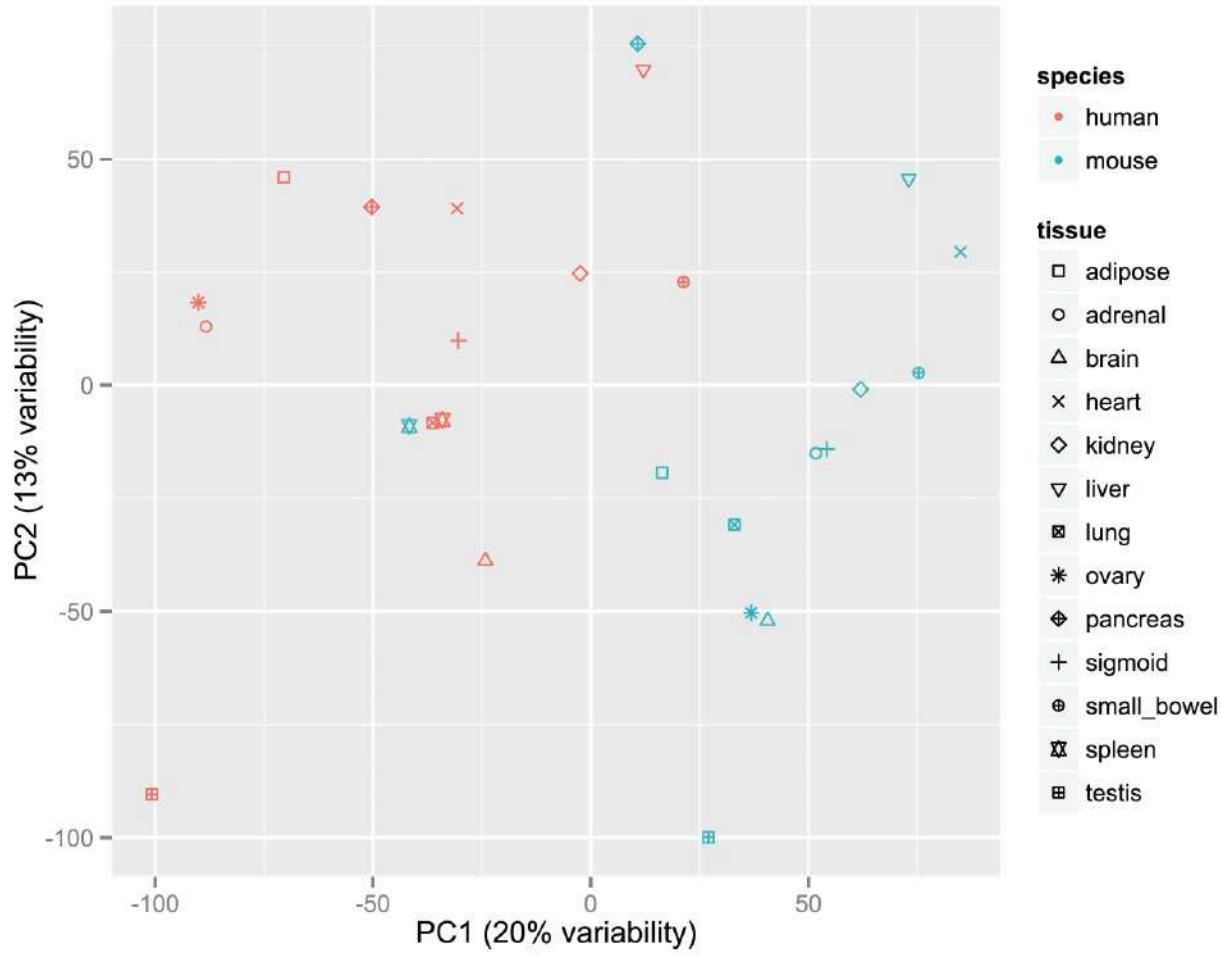
Example of batch effect:



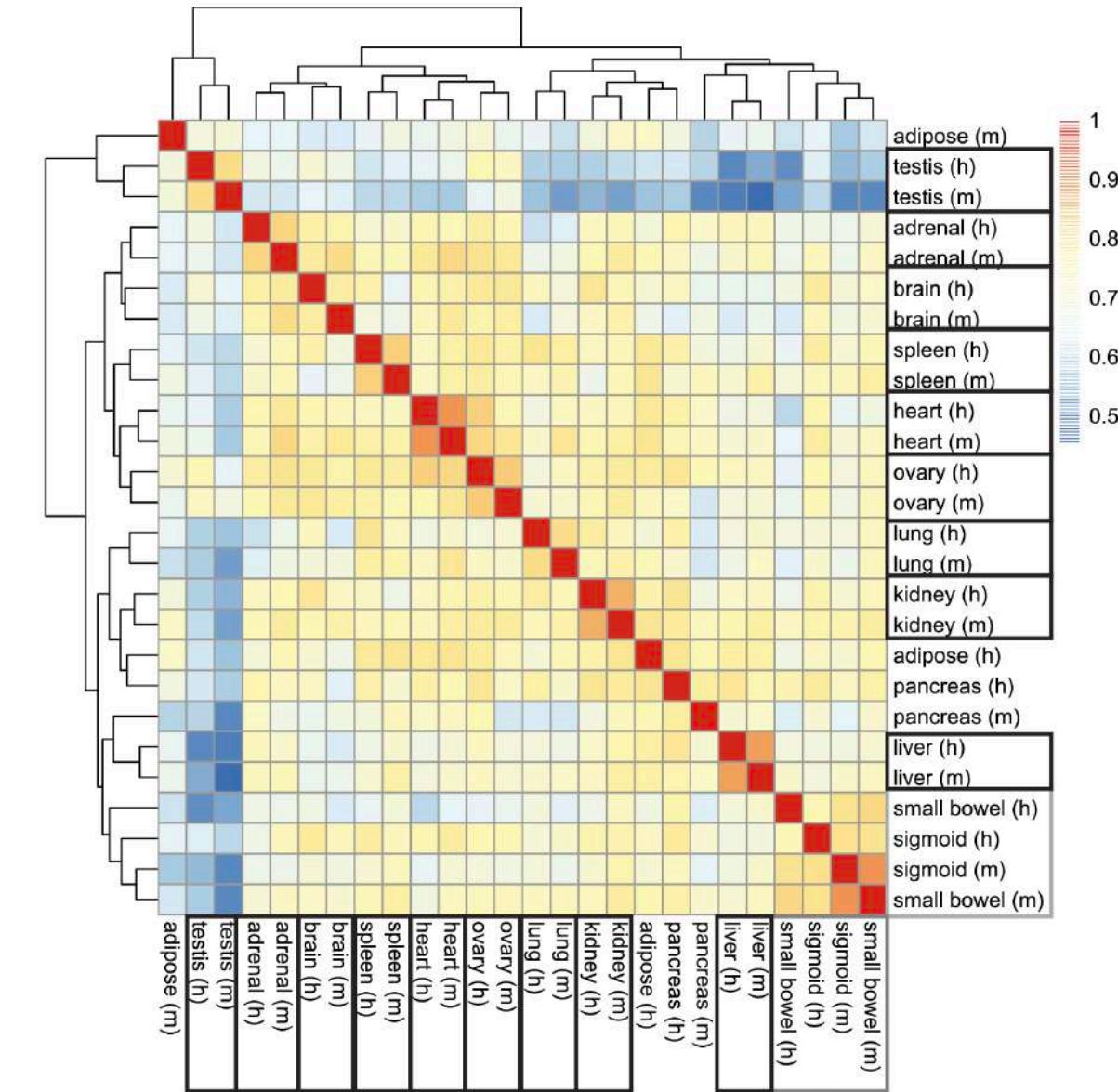
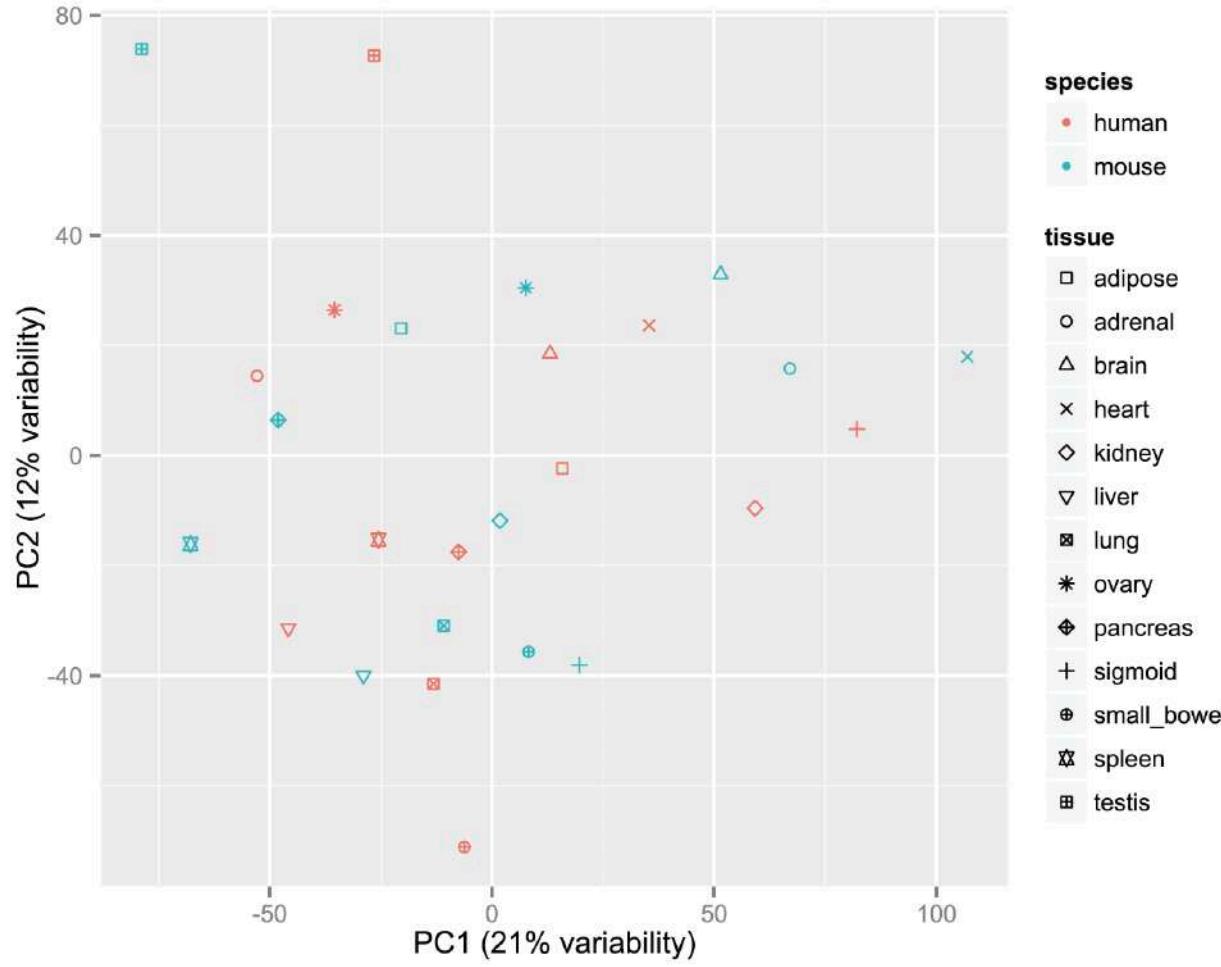
Example of batch effect:

D87PMJN1 (run 253, flow cell D2GUAACXX, lane 7)	D87PMJN1 (run 253, flow cell D2GUAACXX , lane 8)	D4LHBFN1 (run 276, flow cell C2HKJACXX , lane 4)	MONK (run 312, flow cell C2GR3ACXX , lane 6)	HWI-ST373 (run 375, flow cell C3172ACXX , lane 7)
heart	adipose	adipose	heart	brain
kidney	adrenal	adrenal	kidney	pancreas
liver	sigmoid colon	sigmoid colon	liver	brain
small bowel	lung	lung	small bowel	spleen
spleen	ovary	ovary	testis	● Human
testis		pancreas		● Mouse

Recapitulating the patterns reported by the mouse ENCODE papers

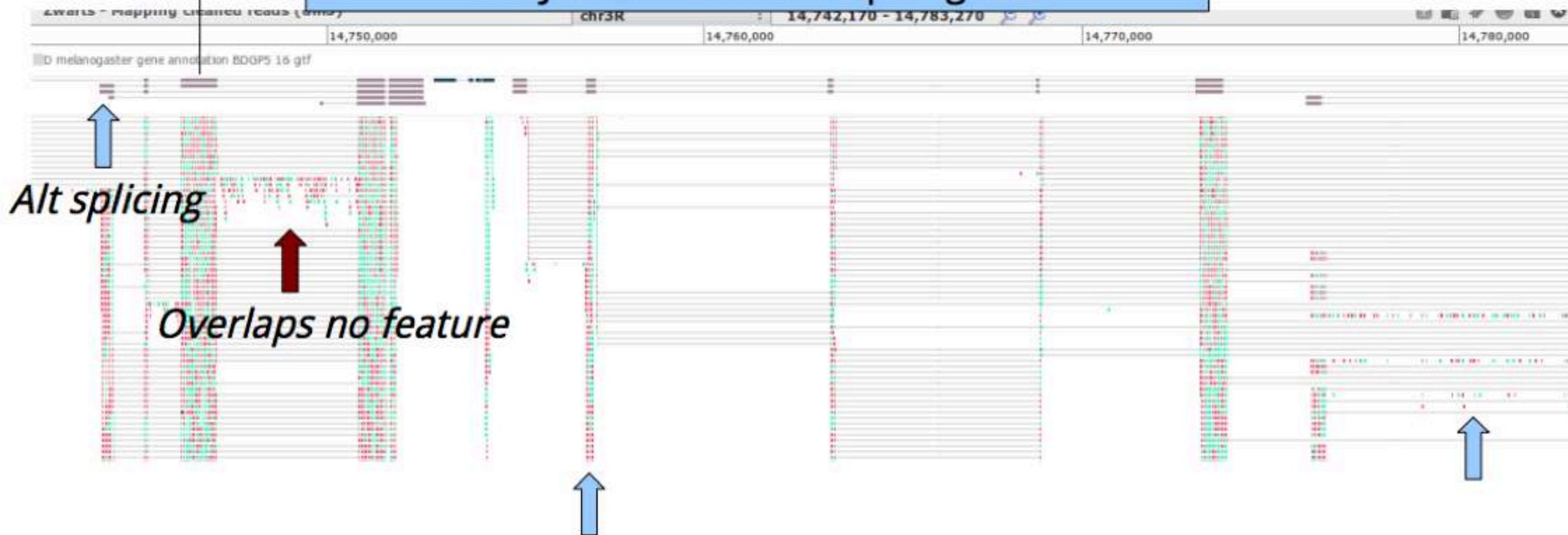


Clustering of data once batch effects are accounted for



Once you have mappings, you can start counting

'Exons' are the type of *features* used here.
They are summarized per 'gene'



Concept:

GeneA = exon 1 + exon 2 + exon 3 + exon 4 = 215 reads

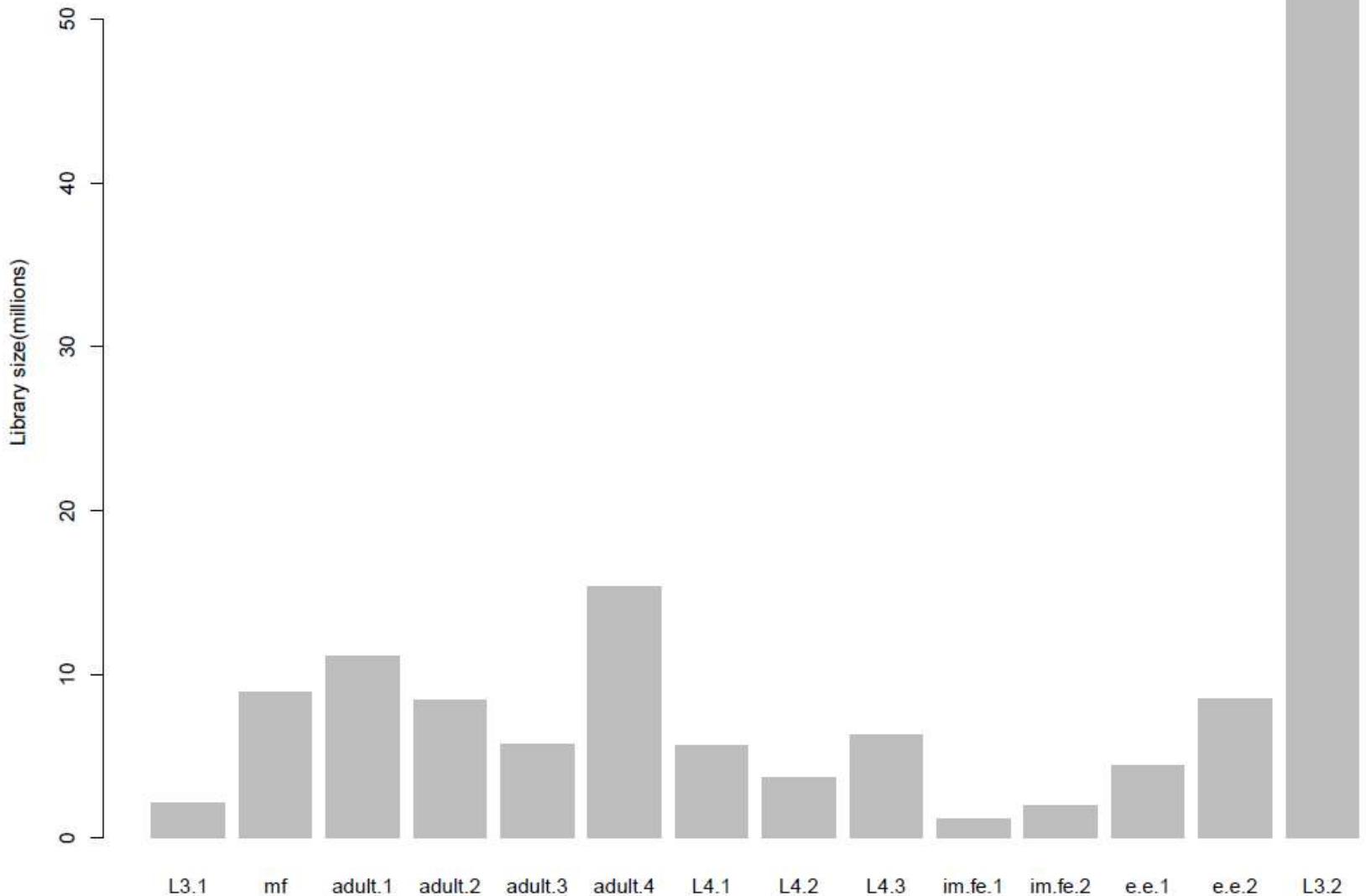
GeneB = exon 1 + exon 2 + exon 3 = 180 reads

This is the bit we care about!

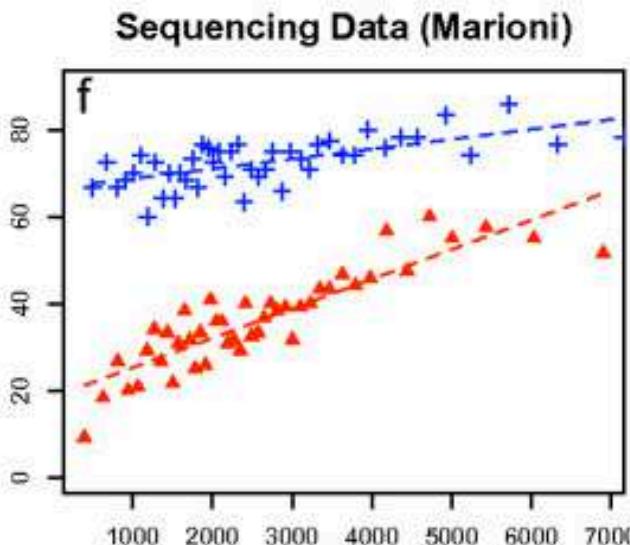
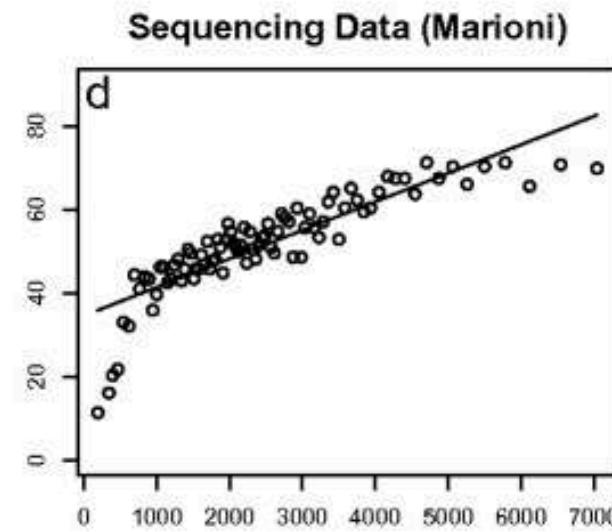
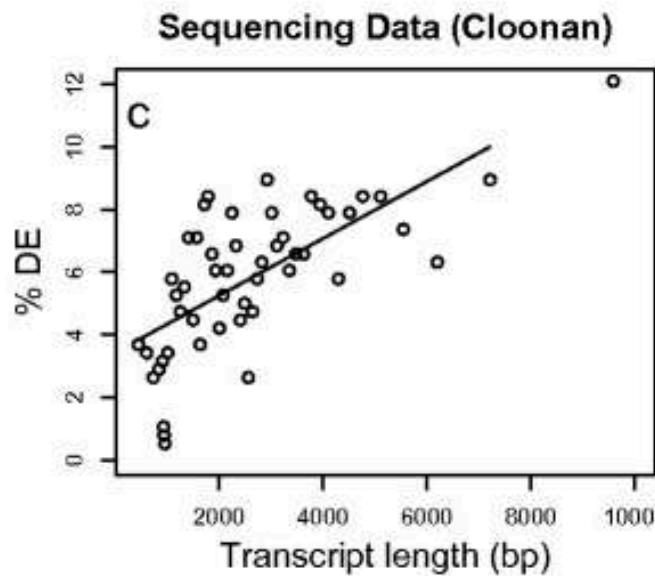
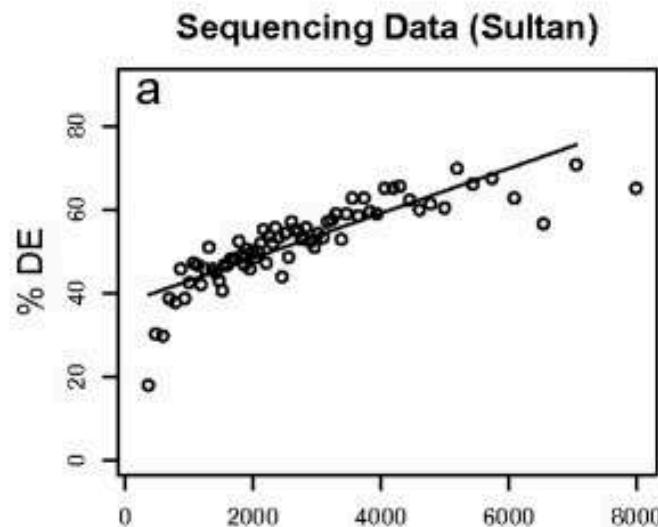


Counts of the gene depends on **expression** ,transcript length
,sequencing depth and simply chance

Higher sequencing depth equals more counts



Counts are proportional to the transcript length x mRNA expression level



33% of highest expressed genes
33% of lowest expressed genes

Normalization: different goals

- **R/FPKM:** (Mortazavi et al. 2008)
 - **Correct for:** differences in sequencing depth and transcript length
 - **Aiming to:** compare a gene across samples and diff genes within sample
- **TMM:** (Robinson and Oshlack 2010)
 - **Correct for:** differences in transcript pool composition; extreme outliers
 - **Aiming to:** provide better across-sample comparability
- **TPM:** (Li et al 2010, Wagner et al 2012)
 - **Correct for:** transcript length distribution in RNA pool
 - **Aiming to:** provide better across-sample comparability
- **Limma voom (logCPM):** (Law et al 2013)
 - **Aiming to:** stabilize variance; remove dependence of variance on the mean

Optimal Scaling of Digital Transcriptomes

Gustavo Glusman , Juan Caballero, Max Robinson, Burak Kutlu, Leroy Hood

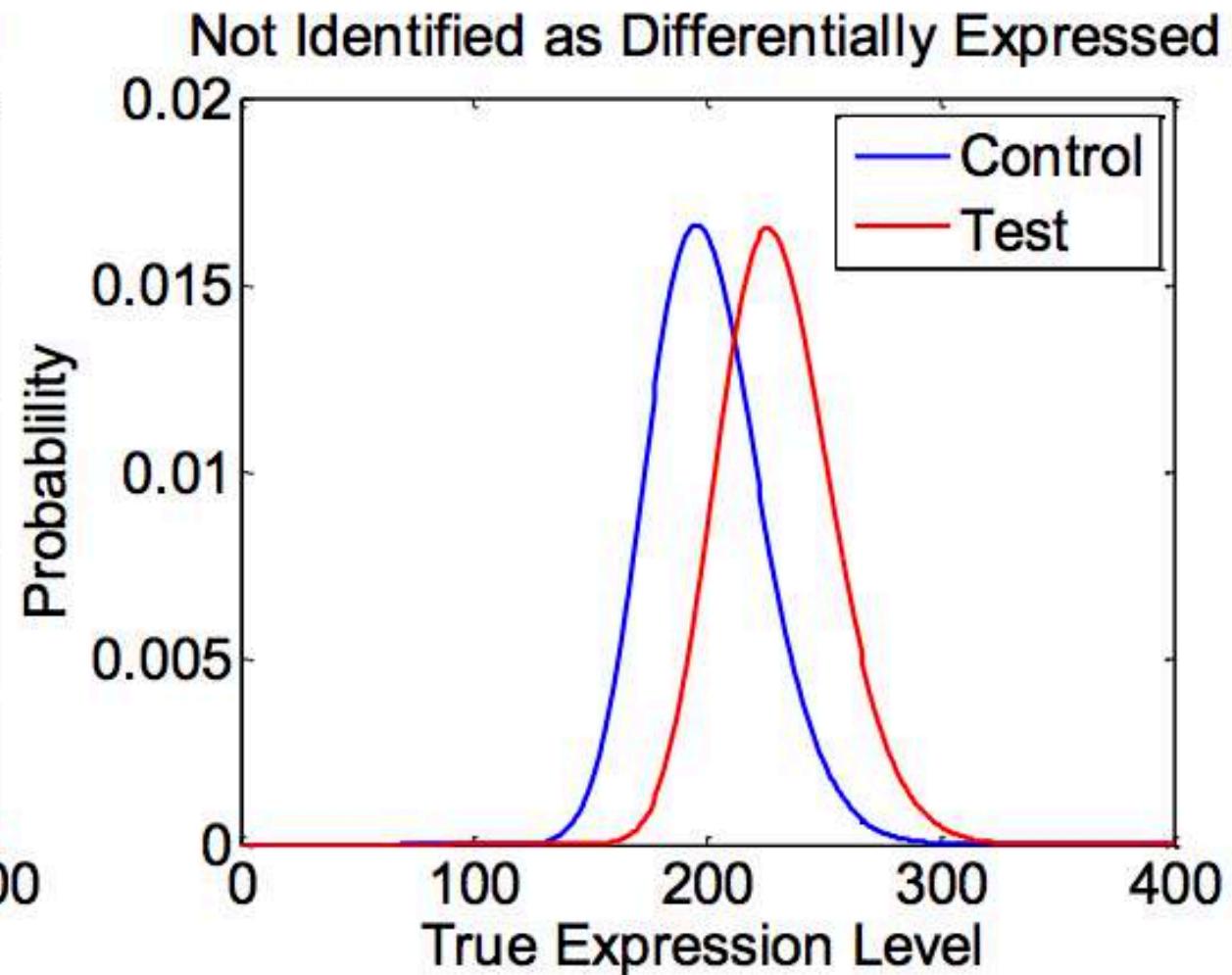
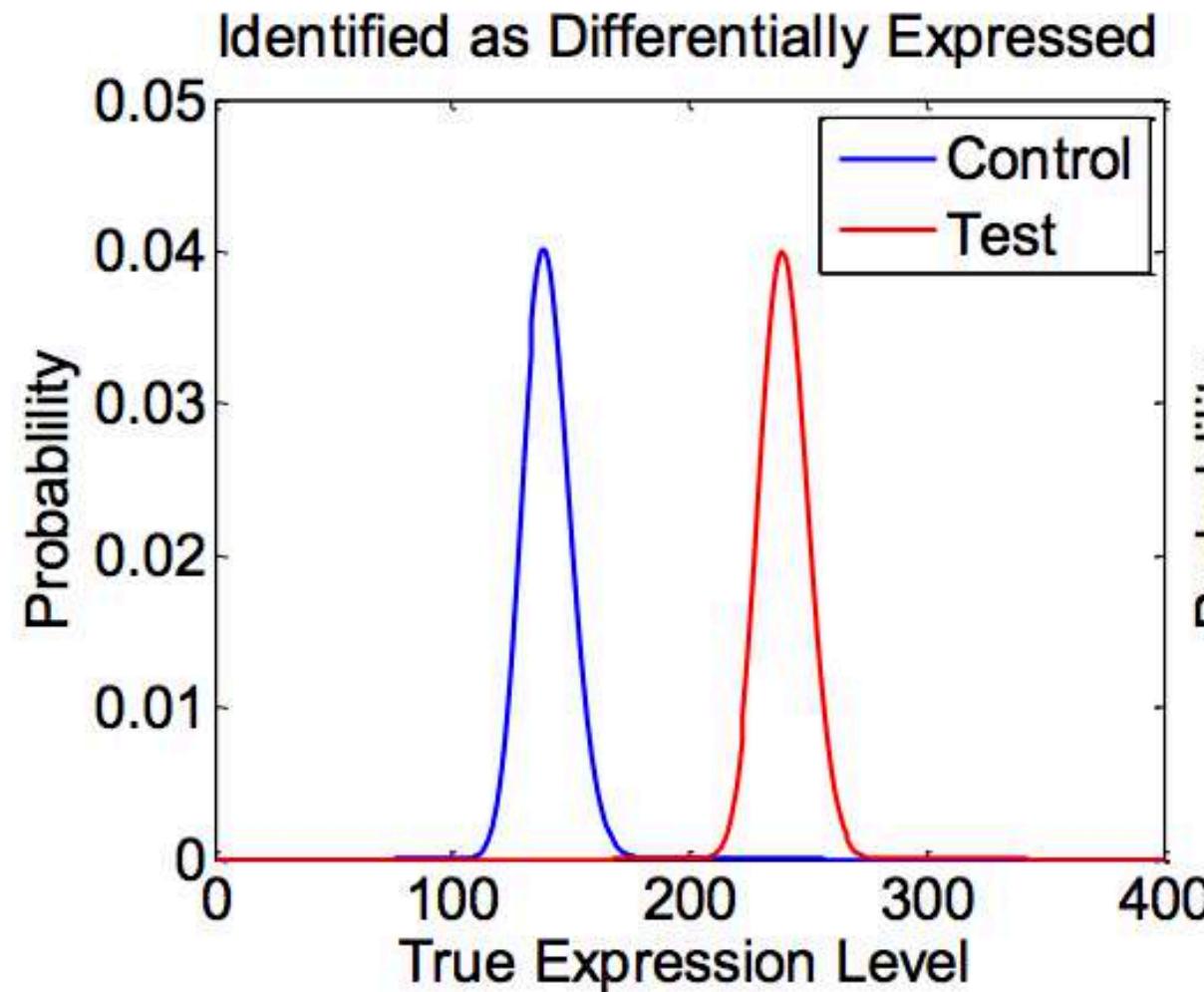
Published: Nov 06, 2013 • DOI: 10.1371/journal.pone.0077885

But how do you know your count = 2 is really 2?

- Differentially expressed genes = counts of genes change between conditions **more systematically** than expected by chance
- Need **biological and technical replicates** to detect differential expression

Differential expression

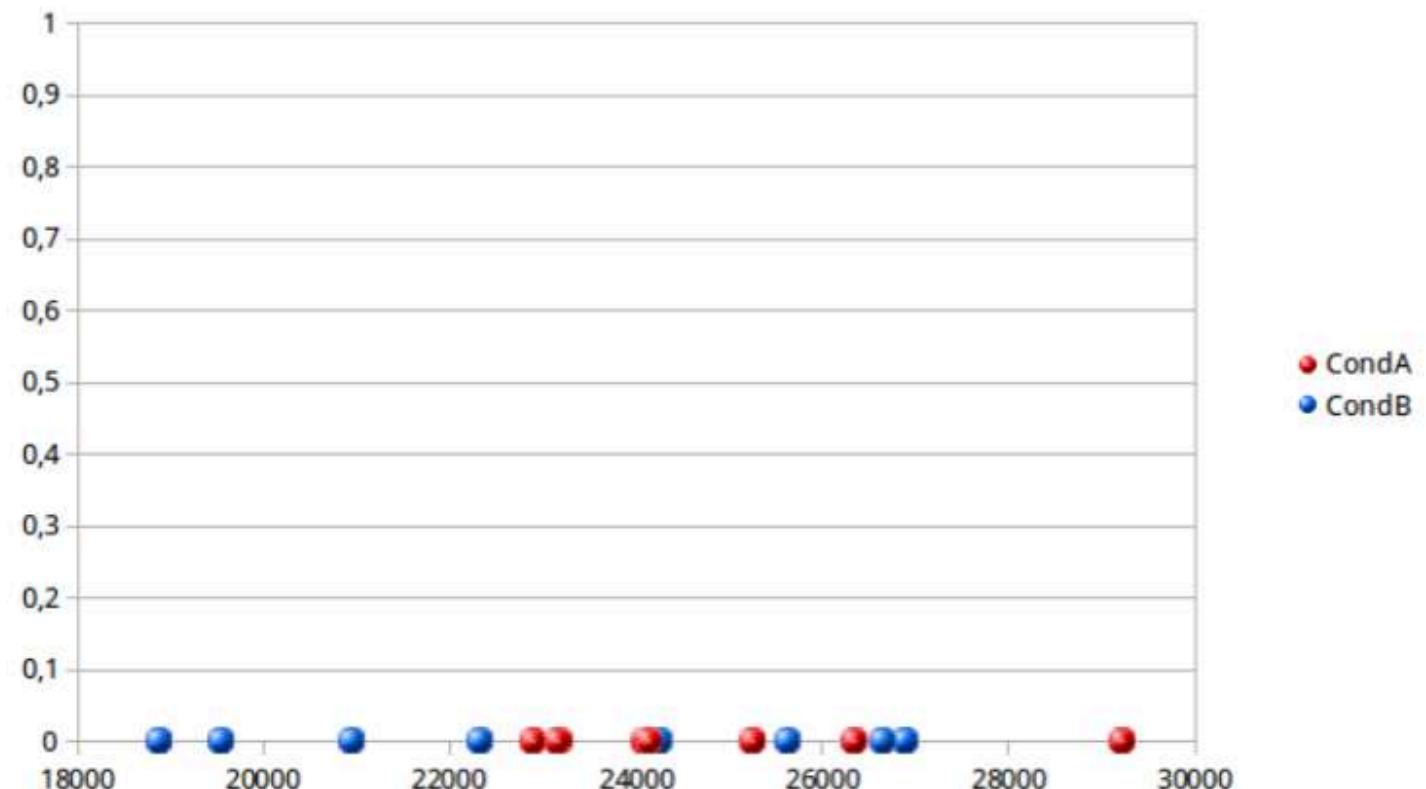
Fitting a distribution for every gene for DE



Scenario

gene_id CAF0006876

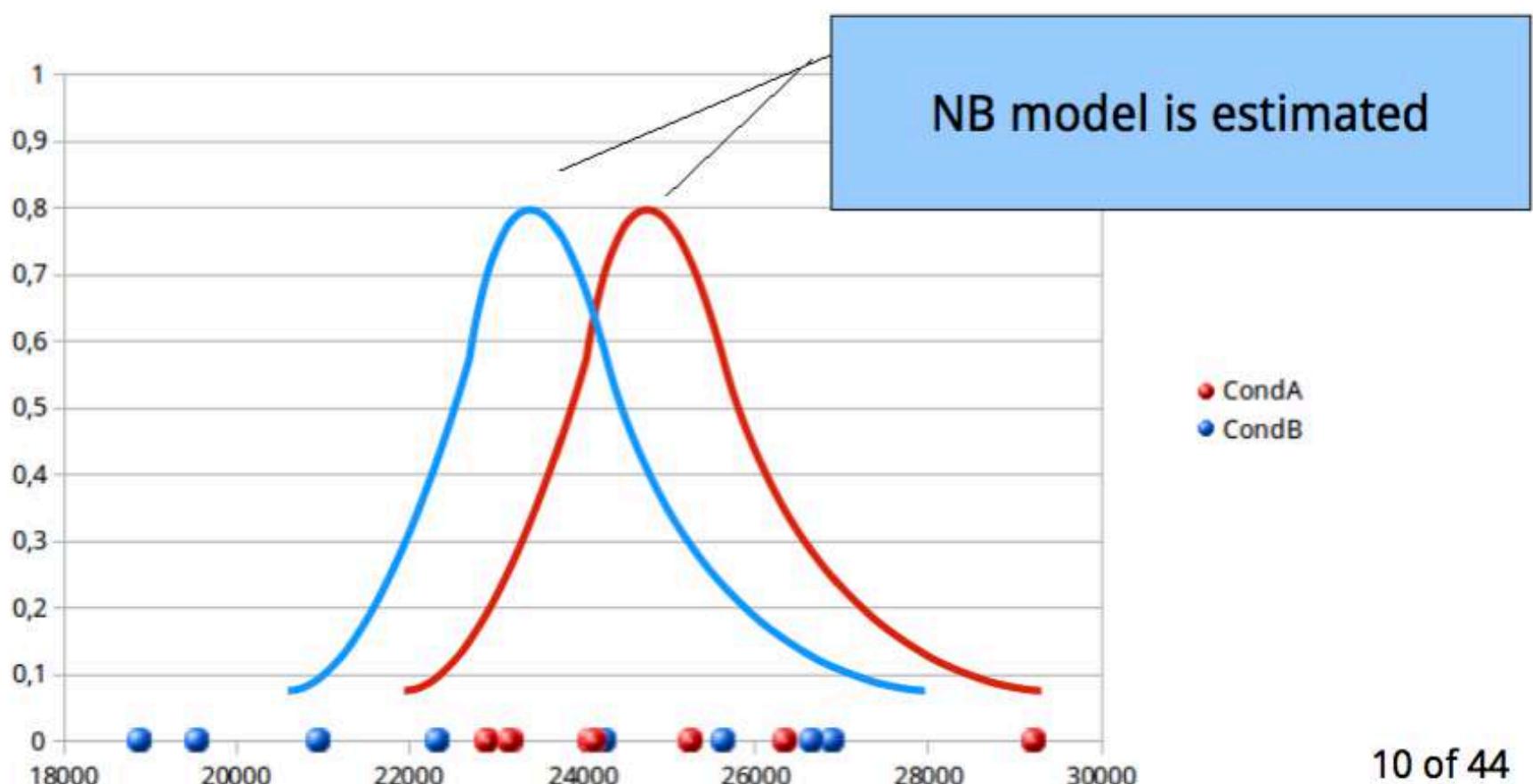
	sample1 23171	sample2 22903	sample3 29227	sample4 24072	sample5 23151	sample6 26336	sample7 25252	sample8 24122
Condition A								
Condition B	Sample9 19527	sample10 26898	sample11 18880	sample12 24237	sample13 26640	sample14 22315	sample15 20952	sample16 25629



Scenario

gene_id CAF0006876

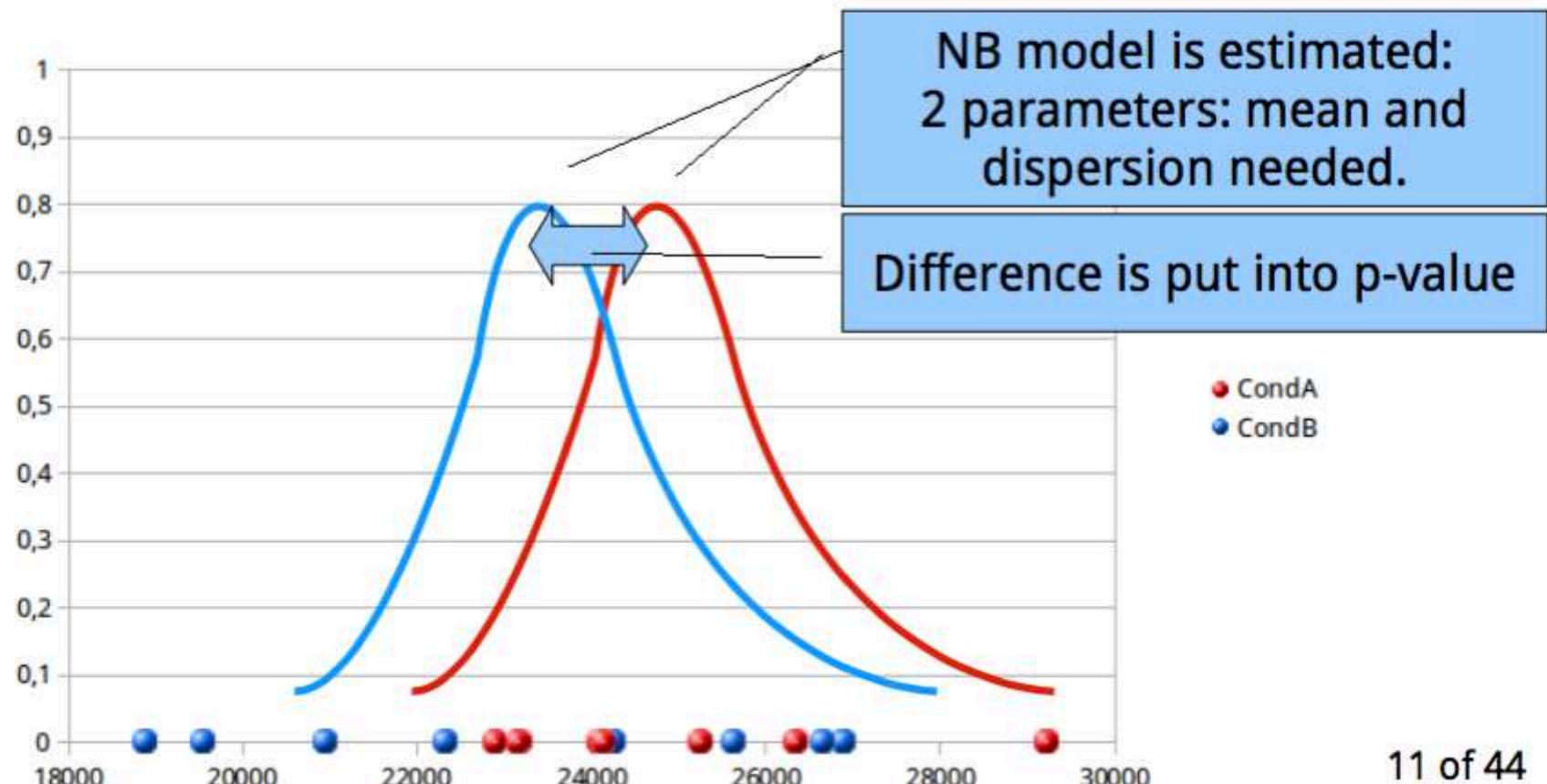
	sample1 23171	sample2 22903	sample3 29227	sample4 24072	sample5 23151	sample6 26336	sample7 25252	sample8 24122
Condition A								
Condition B	Sample9 19527	sample10 26898	sample11 18880	sample12 24237	sample13 26640	sample14 22315	sample15 20952	sample16 25629



Scenario

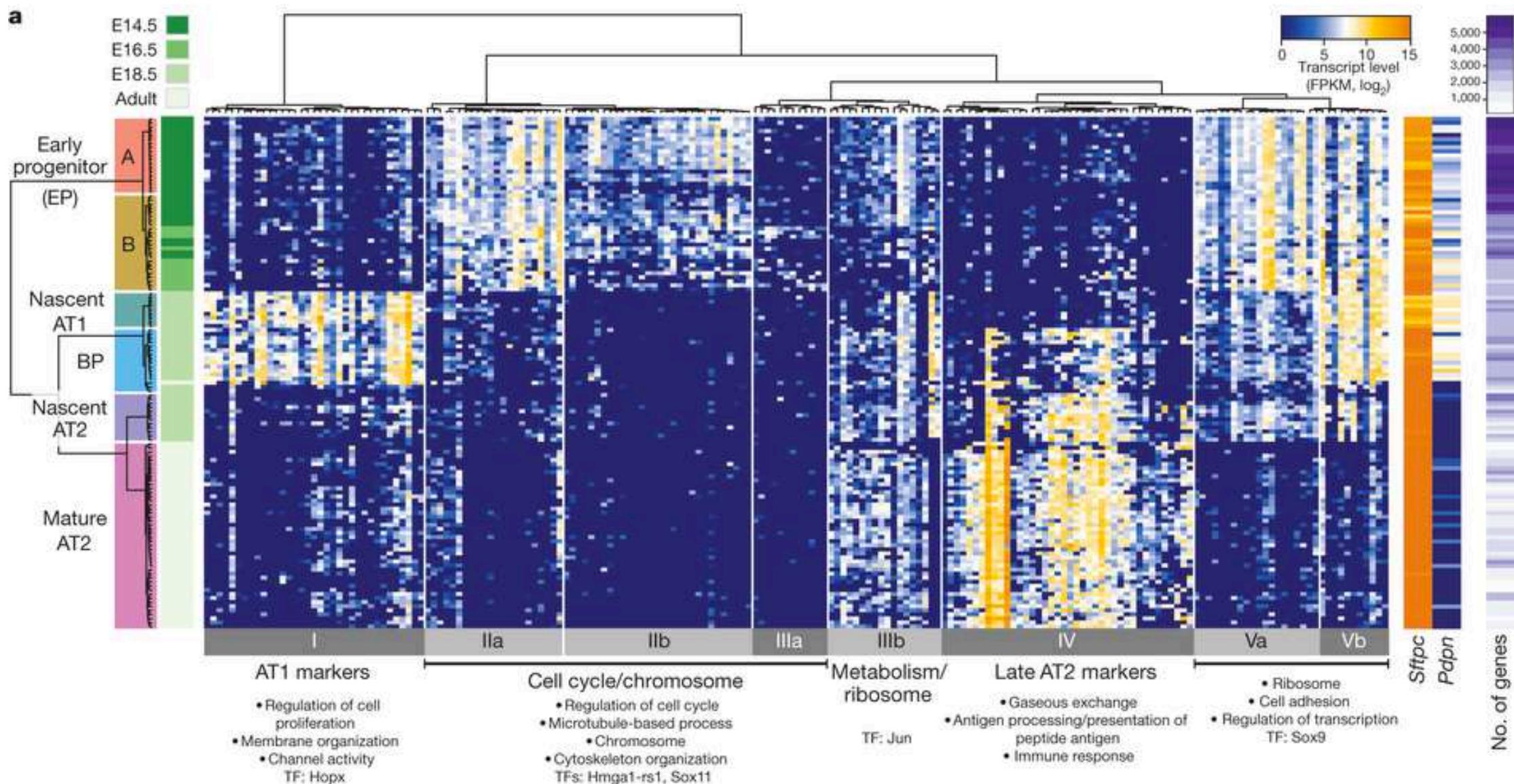
gene_id CAF0006876

	sample1	sample2	sample3	sample4	sample5	sample6	sample7	sample8
Condition A	23171	22903	29227	24072	23151	26336	25252	24122
Condition B	Sample9 19527	sample10 26898	sample11 18880	sample12 24237	sample13 26640	sample14 22315	sample15 20952	sample16 25629



Once you have set of differentially expressed genes

Summarization visualizing the expression data through heatmap ; Classification using Gene Ontology terms and metabolic annotations



Amplicon / Metagenomics: An Intro



A clinician's guide to microbiome analysis

Marcus J. Claesson^{1,2}, Adam G. Clooney^{1-3*} and Paul W. O'Toole^{1,2}*

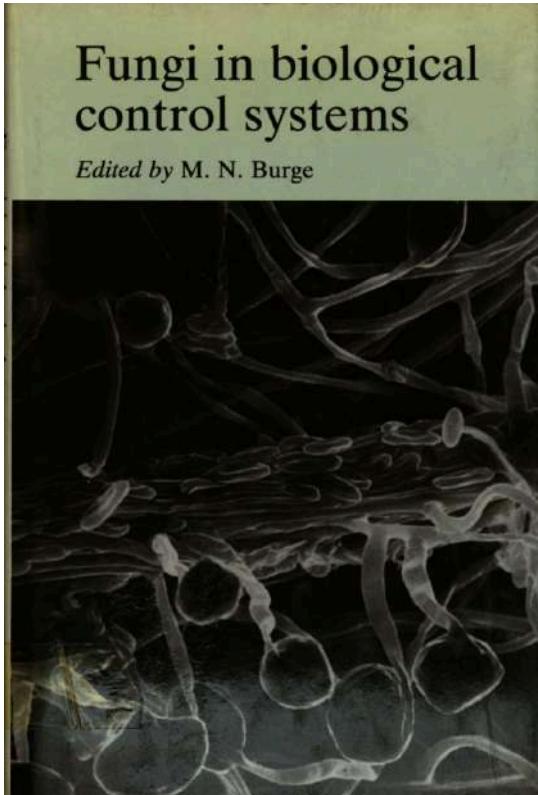
Abstract | Microbiome analysis involves determining the composition and function of a community of microorganisms in a particular location. For the gastroenterologist, this technology opens up a rapidly evolving set of challenges and opportunities for generating novel insights into the health of patients on the basis of microbiota characterizations from intestinal, hepatic or extraintestinal samples. Alterations in gut microbiota composition correlate with intestinal and extraintestinal disease and, although only a few mechanisms are known, the microbiota are still an attractive target for developing biomarkers for disease detection and management as well as potential therapeutic applications. In this Review, we summarize the major decision points confronting new entrants to the field or for those designing new projects in microbiome research. We provide recommendations based on current technology options and our experience of sequencing platform choices. We also offer perspectives on future applications of microbiome research, which we hope convey the promise of this technology for clinical applications.

Key points

- Complex communities of microorganisms live on and in the human body, and variations in the composition and function of these communities are increasingly linked to various conditions and diseases
- Although it is not known if microbiome changes are causative or consequential in most pathophysiologies, they might provide biomarkers for disease detection or management
- Microbiome analysis is likely to become a routine component of secondary health care and is emerging as a modifiable environmental risk factor in multifactorial diseases that could be targeted by novel therapeutics
- Technology advancements are leading to a range of powerful methods for microbiome analysis becoming available and affordable for clinical studies
- Judicious choice of sample type and sequencing platform are required to maximize the clinical utility of microbiome data

What is the microbiome?

Fungi in Biological Control Systems (1988)



A convenient ecological framework in which to examine biocontrol systems is that of the microbiome. This may be defined as a characteristic microbial community occupying a reasonably well defined habitat which has distinct physico-chemical properties. The term thus not only refers to the microorganisms involved but also encompasses their theatres of activity. In relation to fungal diseases of crops and their control, major microbiomes are the phylloplane, spermosphere, rhizosphere and rhizoplane, and numerous kinds of plant residues persisting on or in the soil. Mention should also be made of the wood of standing or felled trees as microbiomes where biocontrol of forest diseases using fungi has been achieved. However, in most cases competitive interactions other than mycoparasitism seem to be of greater importance.

<http://microbe.net/2015/04/08/what-does-the-term-microbiome-mean-and-where-did-it-come-from-a-bit-of-a-surprise/>

And then what is the metagenome?

Crosstalk R245

Molecular biological access to the chemistry of unknown soil microbes: a new frontier for natural products

Jo Handelsman¹, Michelle R Rondon¹, Sean F Brady², Jon Clardy² and Robert M Goodman¹

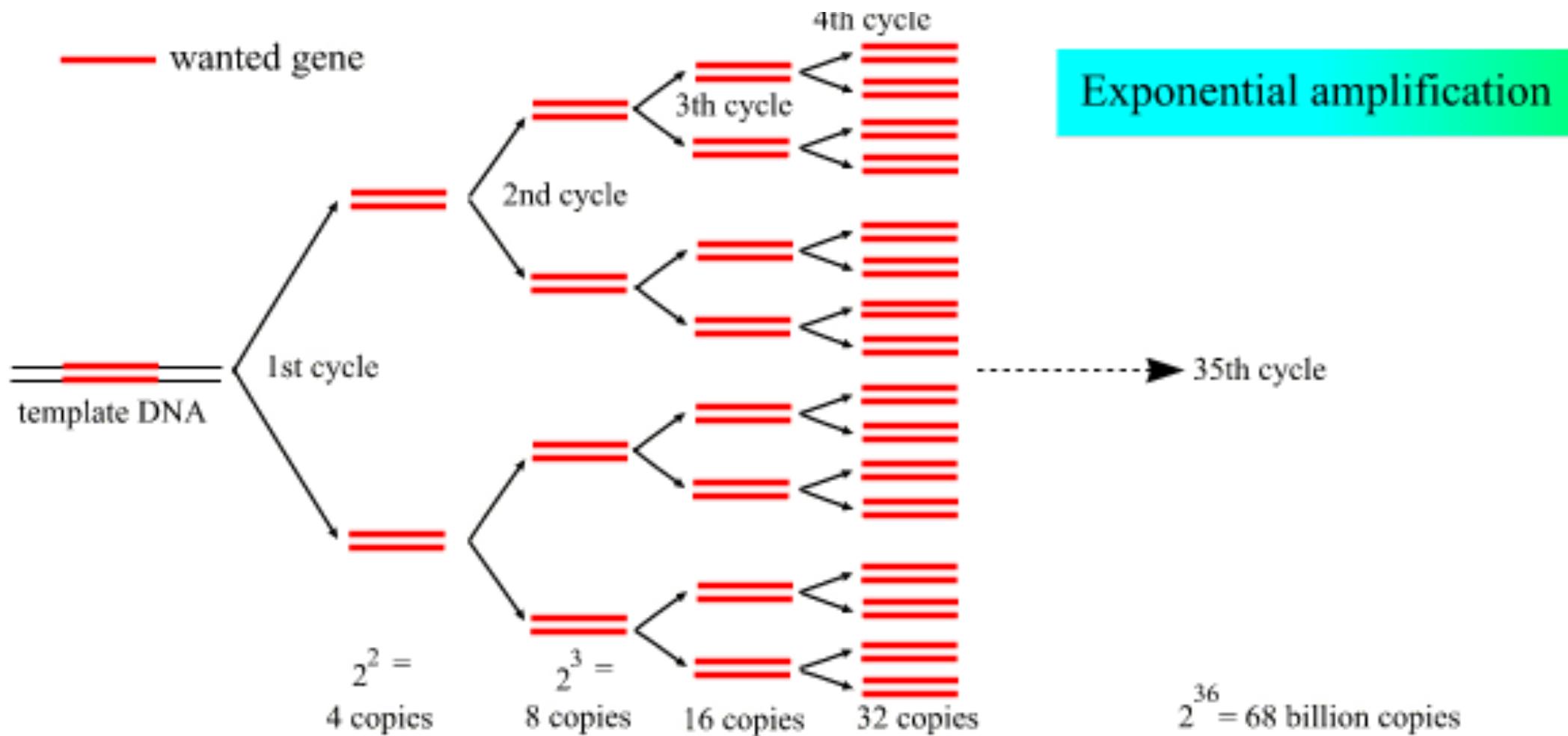


Chemistry & Biology October 1998, 5:R245–249
<http://biomednet.com/electref/10745521005R0245>

... This approach involves directly accessing the genomes of soil organisms that cannot be, or have not been, cultured by isolating their DNA

What is amplicon sequencing?

Anything that requires PCR-based amplification of a specific target gene (locus)



And then what is the metagenome?

OPEN  ACCESS Freely available online

PLOS COMPUTATIONAL BIOLOGY

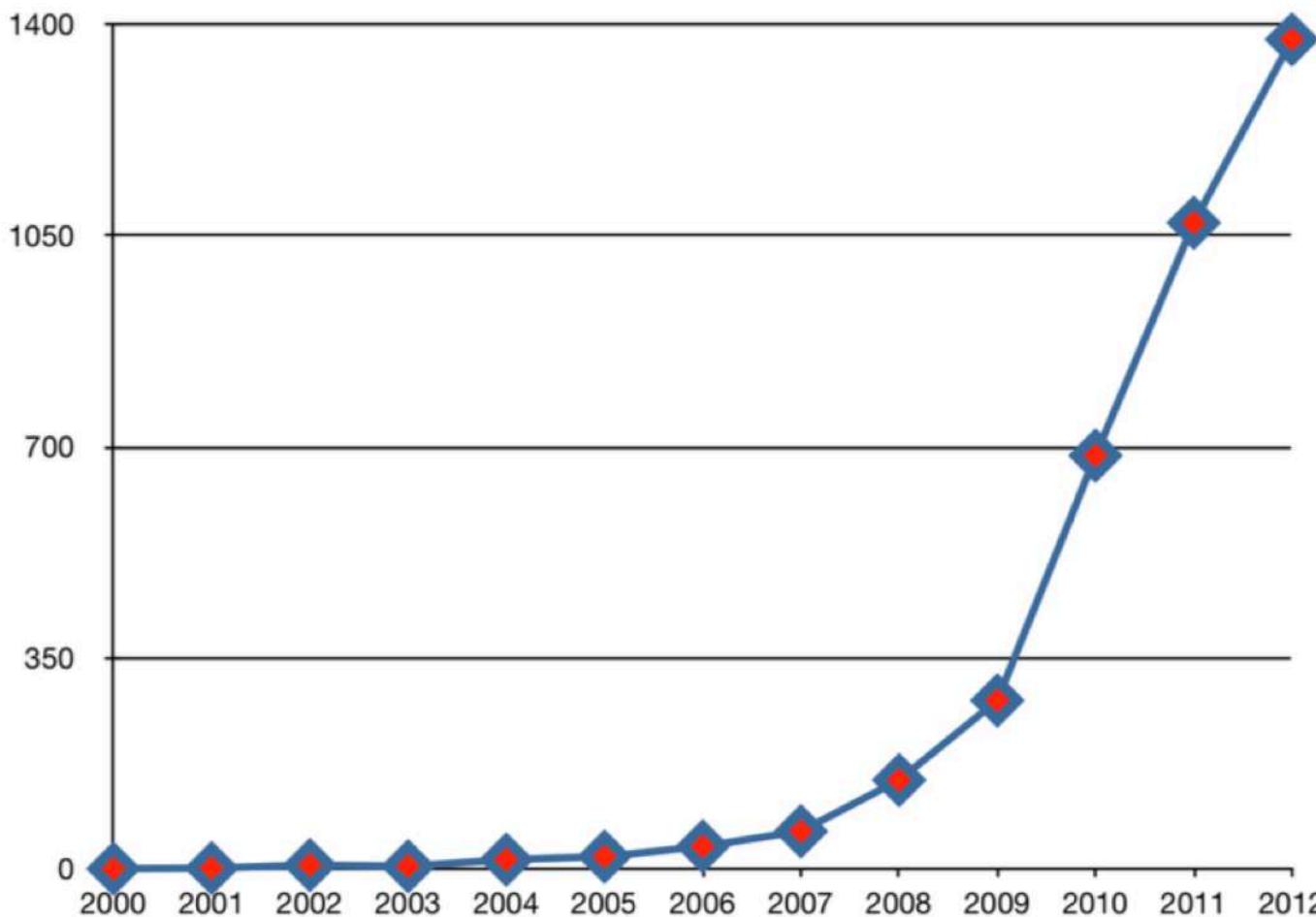
Review

Bioinformatics for Whole-Genome Shotgun Sequencing of Microbial Communities

Kevin Chen*, Lior Pachter*

Metagenomics is the application of modern genomics techniques to the study of communities of microbial organisms directly in their natural environments, bypassing the need for isolation and lab cultivation of individual species. The field has its roots in the culture-independent retrieval of **16S rRNA** genes, pioneered by Pace and colleagues two decades ago.

Pubmed hits for “Microbiome”



Metagenomics \neq Amplicon sequencing

Metagenomics is undergoing a crisis

Please don't make things worse ☺

- Crisis 1
 - **The correlation/causation fallacy.** For example....
 - Patients with type II diabetes have a different gut microbiome compared to healthy patients
 - Does the microbiome cause diabetes?
 - Or do they have a different microbiome because they have diabetes? (therefore different diet)
- Crisis 2
 - A lot of people want to do it, but don't know how
 - Errors, bad experimental design, incorrect conclusions

Basic Purpose

Characteristics of (microbial) community

Who are they?

Where do they come from?

Are their similarities (at what level)

between communities

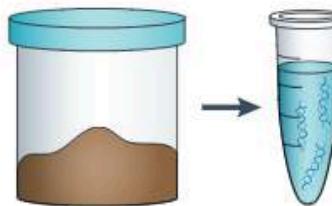
of different conditions

of similar conditions?

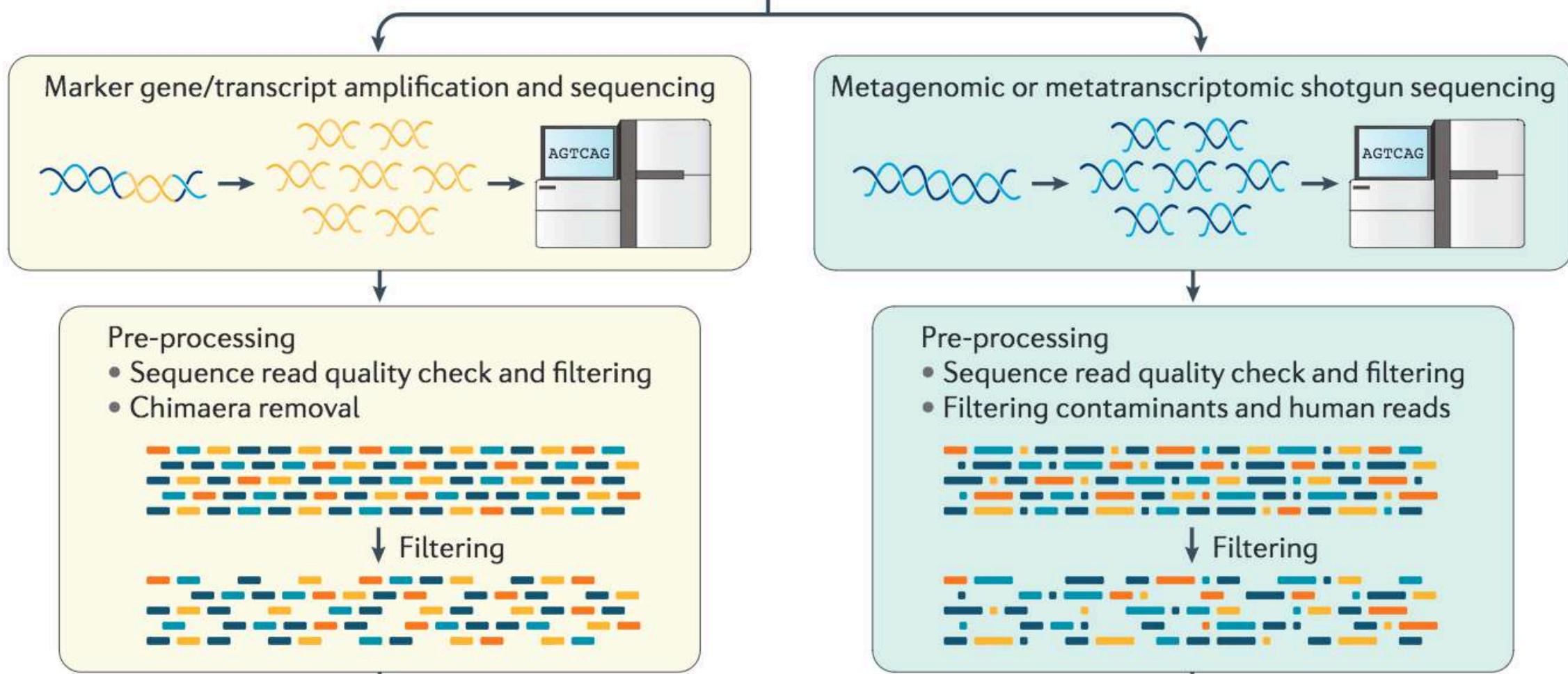
within a community?

What are they doing?

How are they doing?



Study design, sample collection, storage and DNA/cDNA/RNA extraction



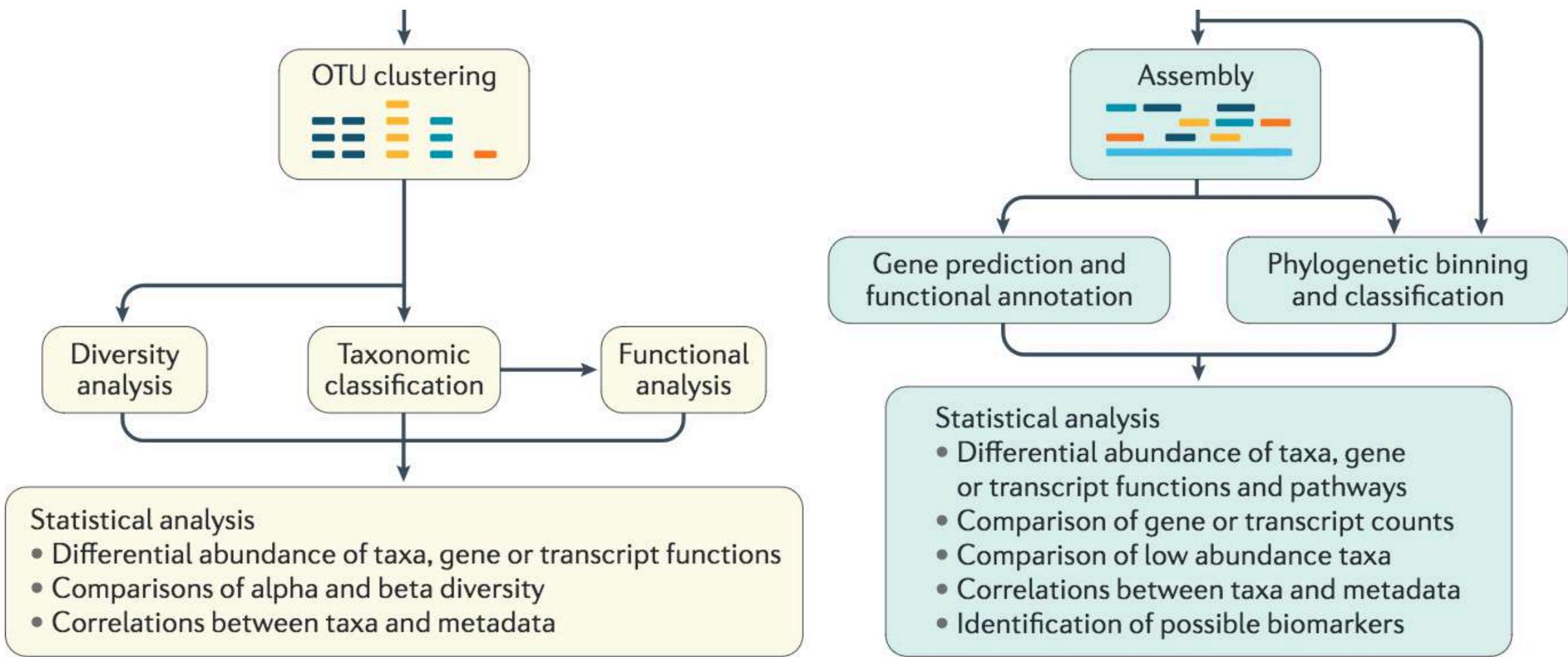


Figure 1 | Flowchart of the major steps involved in bioinformatic analysis of the microbiome. The analysis is divided into two sections depending on the type of sequencing. This schematic describes the basic steps and might vary depending on the aim of the analysis. OTU, operational taxonomic unit.

Applications

What have metagenomics been used for?

Exploration and categorisation



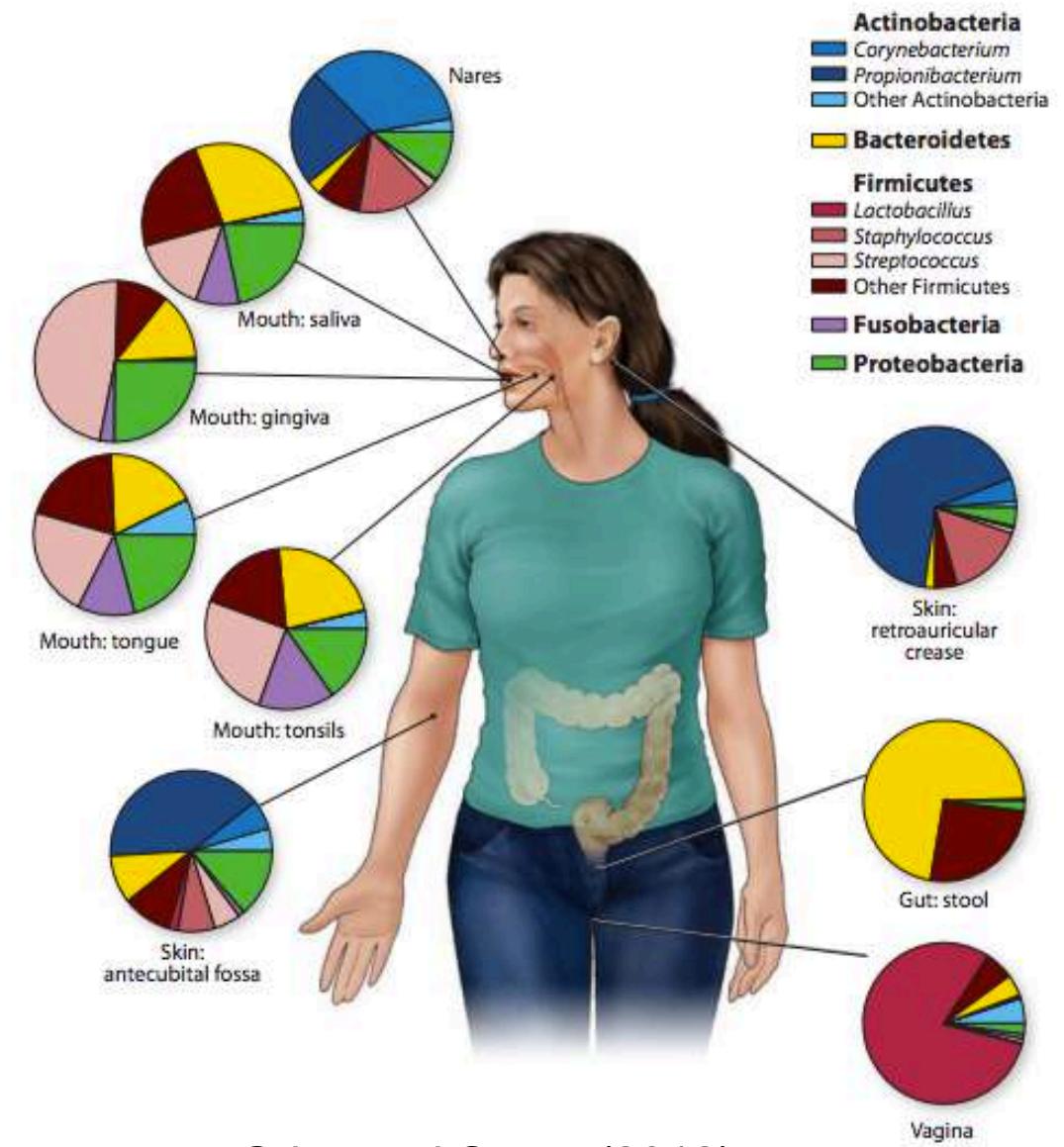
Rusch et al., 2007 Plos Biology

- 6.3 Gbp of sequence (2x Human genomes, 2000 x Bacterial genomes)
- Most sequences were novel compared to the databases



Qin et al., 2010 Nature

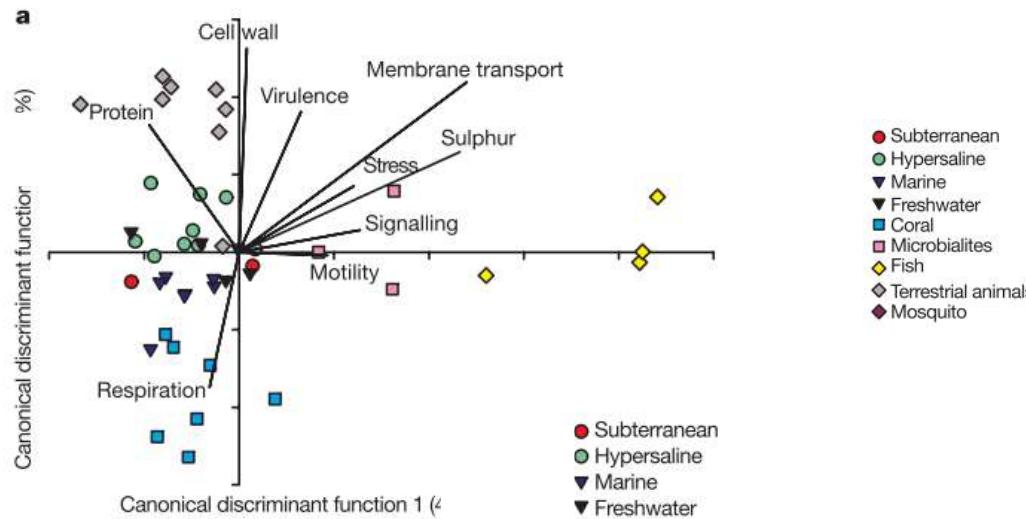
- 127 Human gut metagenomes
- 600 Gbp sequence (200 x Human genomes)
- 3.3 million genes identified
- Minimal gut metagenome defined



Grice and Segre (2012)

What have metagenomics been used for?

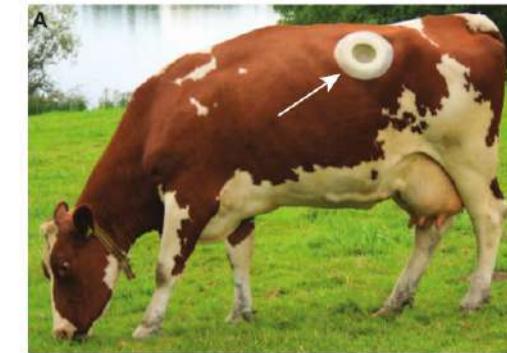
Comparative



Dinsdale *et al.*, 2008 **Nature**

- A characteristic microbial fingerprint for each of the nine different ecosystem types

Specific functions

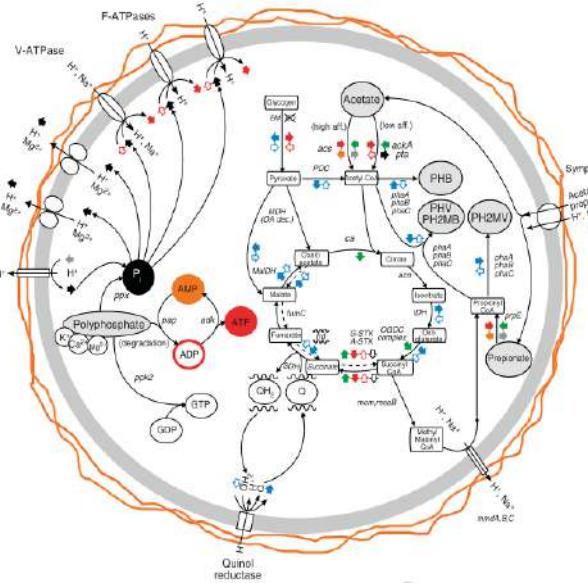


Hess *et al.*, 2011 **Science**

- Identified 27,755 putative carbohydrate-active genes from a cow rumen metagenome
- Expressed 90 candidates of which 57% had enzymatic activity against cellulosic substrates

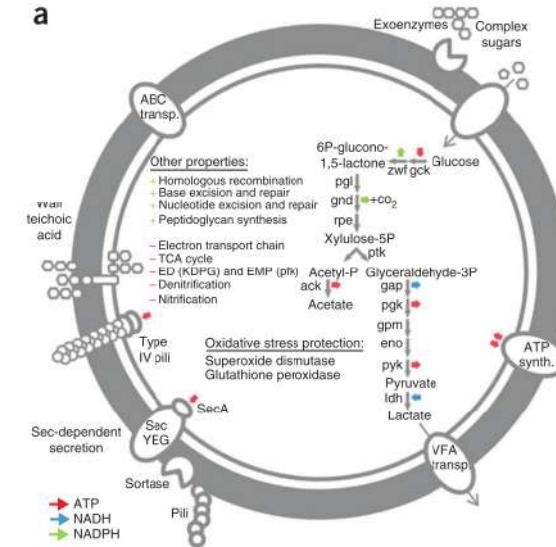
What have metagenomics been used for?

Extracting genomes



Garcia Martin *et al.*, 2006 **Nat. Biotechnol.** Albertsen *et al.*, 2013 **Nat. Biotechnol.**

- Genome extraction from low complexity metagenome
 - *Candidatus Accumulibacter phosphatis*
 - The first genome of a polyphosphate accumulating organism (PAO) with a major role en enhanced biological phosphorus removal



- Genome extraction of low abundant species (< 0.1%) from metagenomes
 - First complete TM7 genome
 - Access to genomes of the "uncultured majority"

Concept: OTU (Operational Taxonomic Unit)

OTU for Ecology

Operational Taxonomic Unit: a grouping of similar sequences that can be treated as a single “species”

Strengths

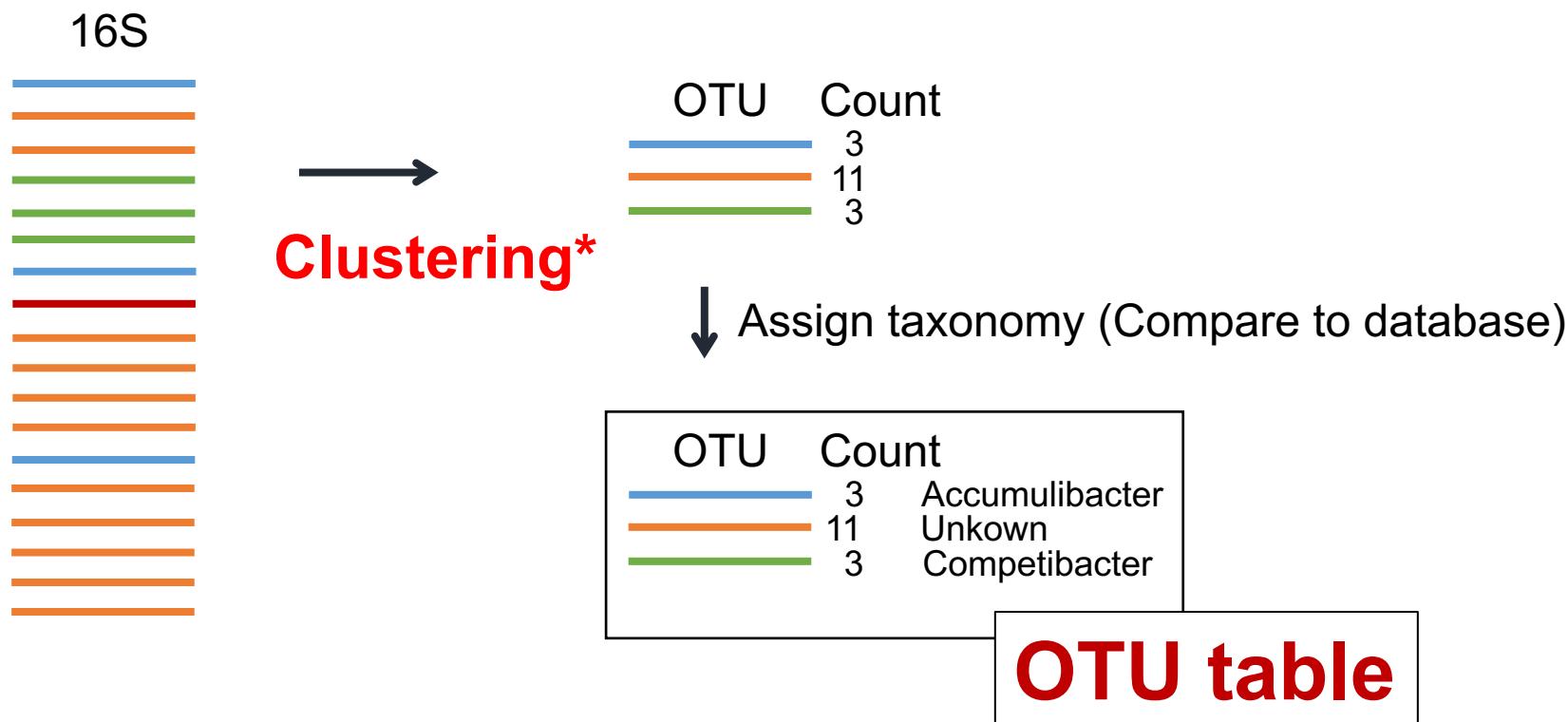
- Conceptually simple
- Mask effect of poor quality data
 - Sequencing error
 - in vitro recombination

Weaknesses

- Limited resolution
- Logically inconsistent definition

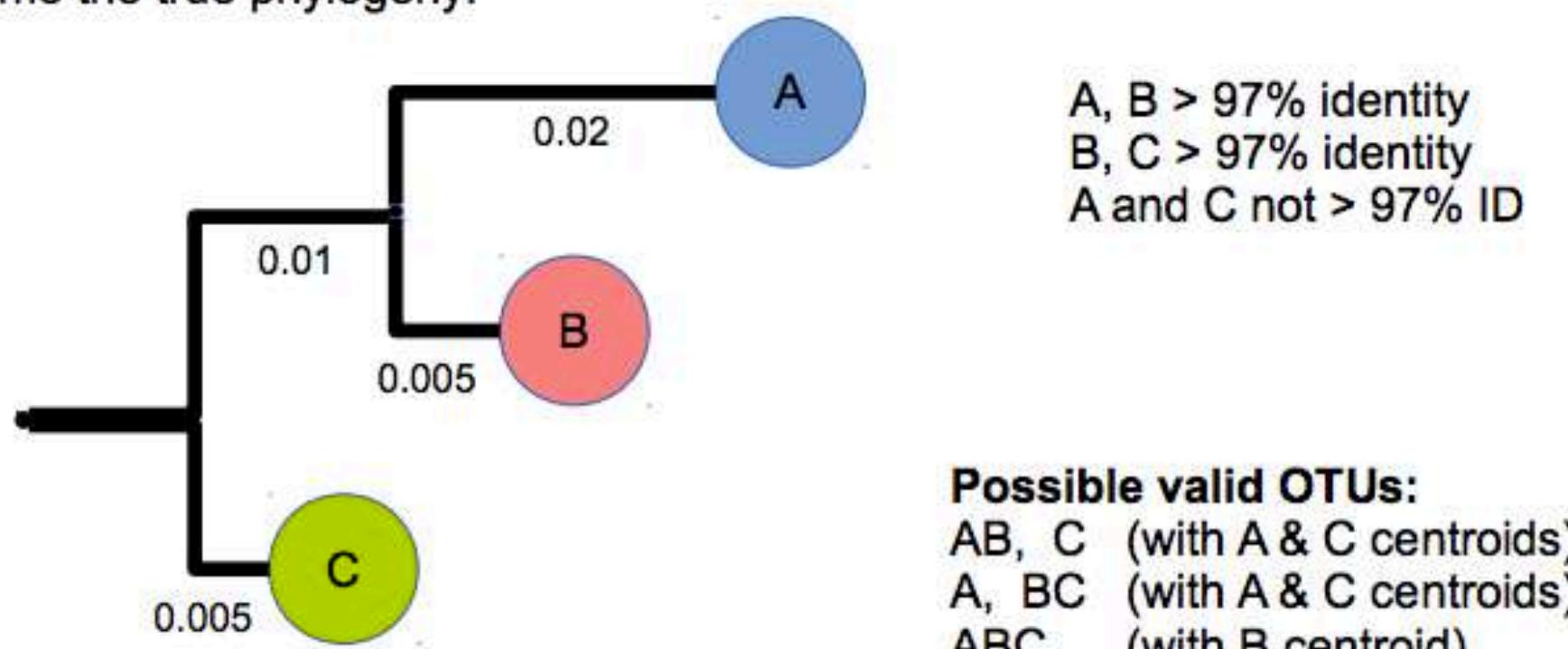
Assign OTU

- Cluster by their similarity to other sequences in the sample (operations taxonomic units → OTU)
- 95% genus level, **97% species level**, 99% strain level



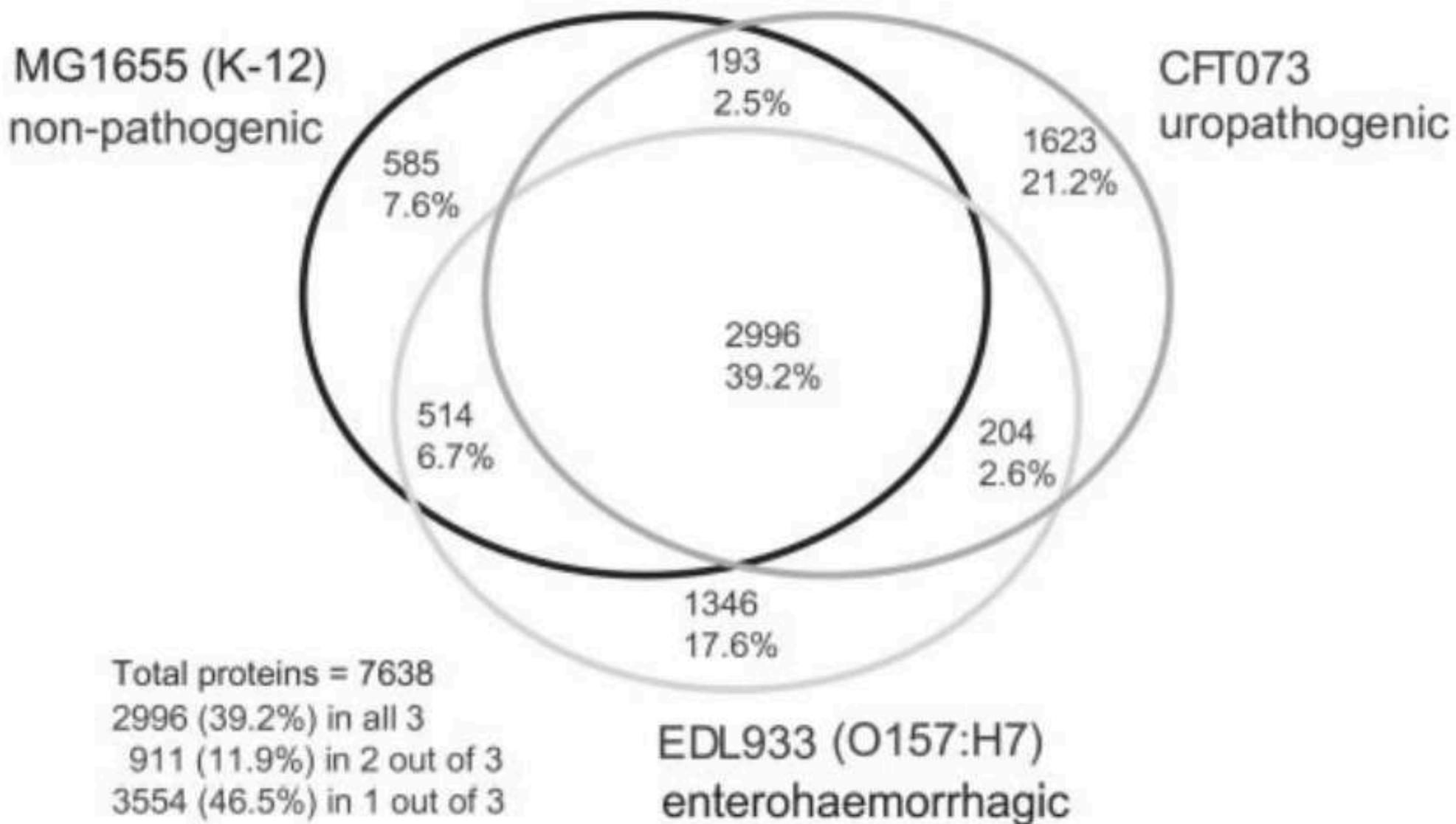
Logical inconsistency: OTUs at 97% ID

Assume the true phylogeny:

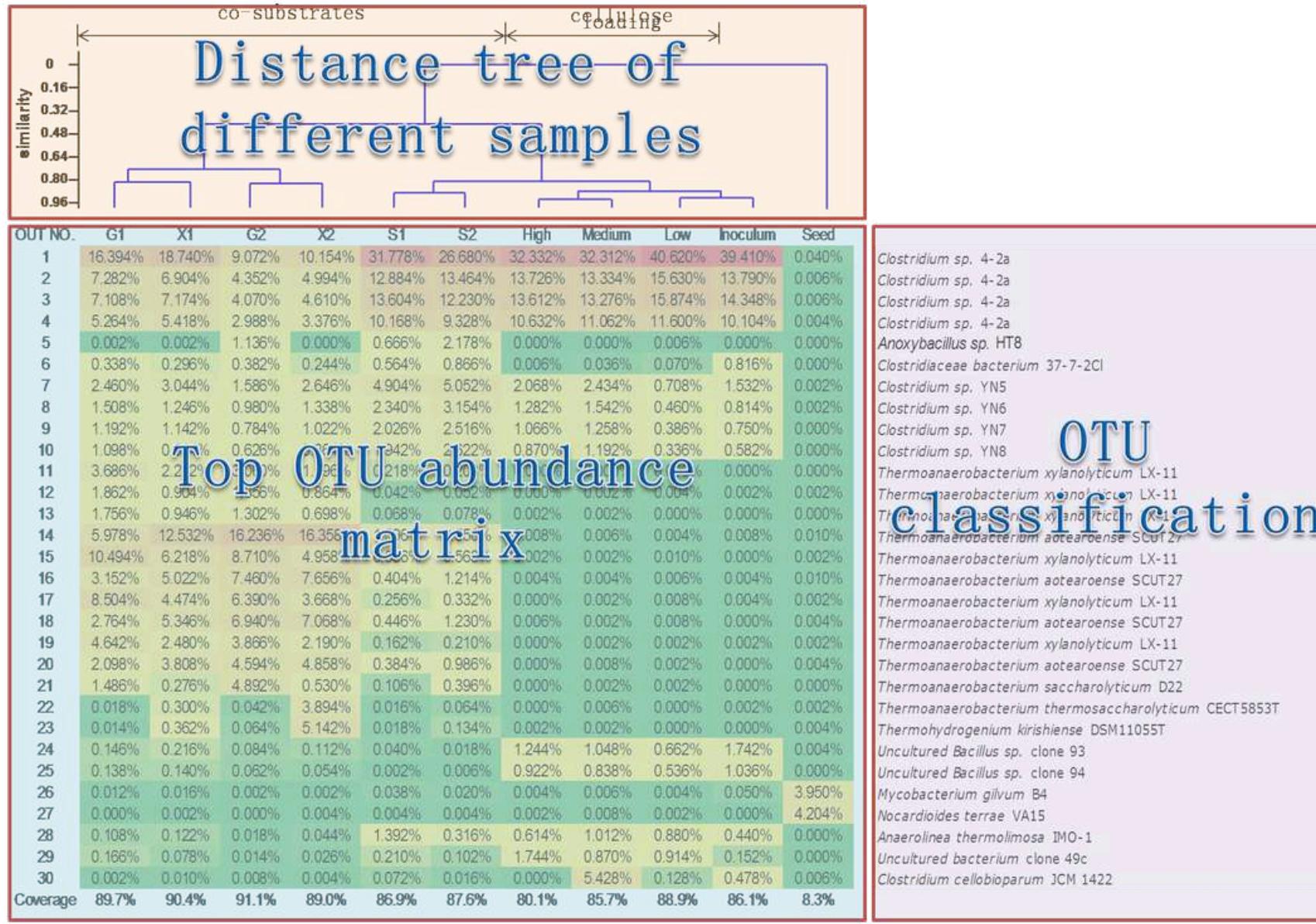


OTU pipelines will arbitrarily pick one of the three solutions.
Is this actually a problem??

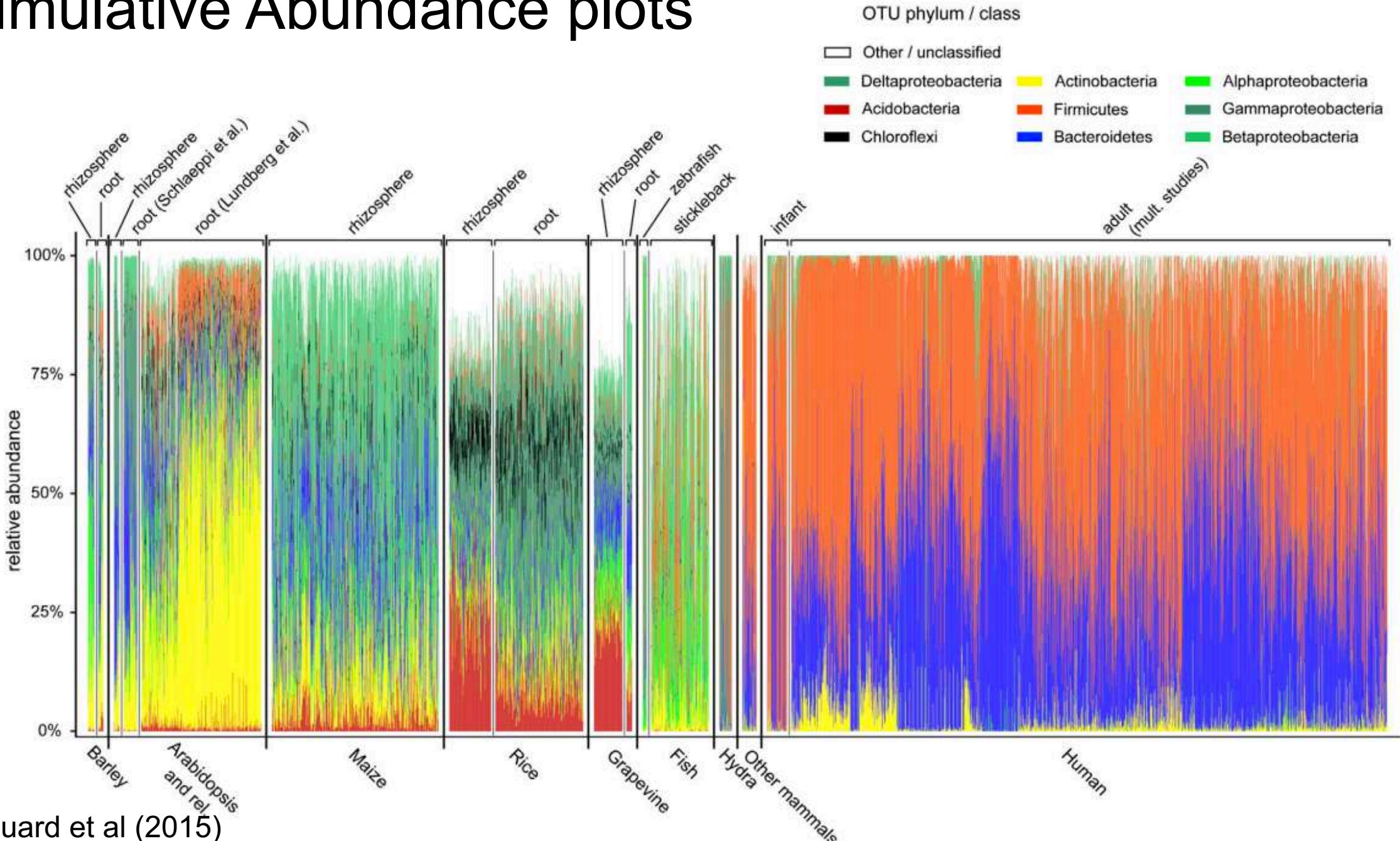
Same species (16S): Different genomes



Tree way plot with top OTUs abundance and classification

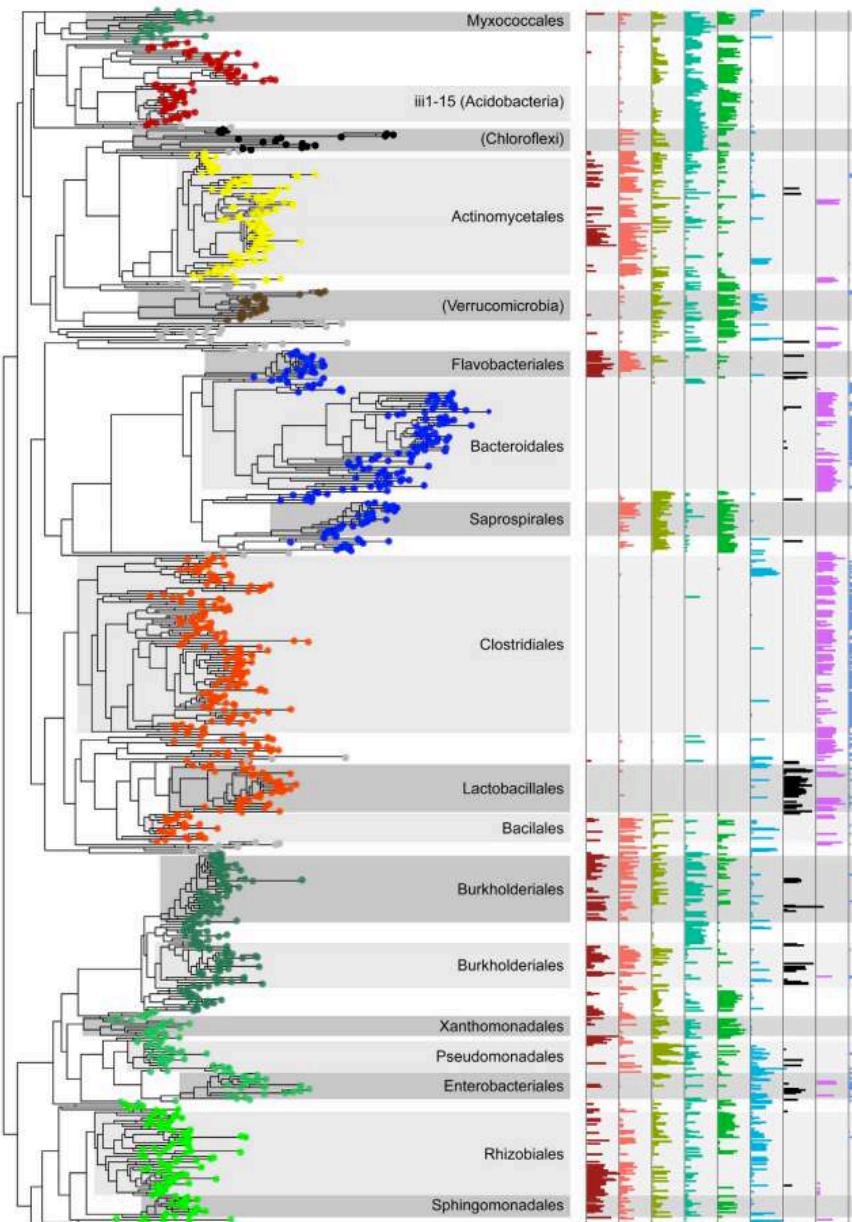


Cumulative Abundance plots



Phylogenetic Analysis of OTU abundances

Relationship between OTUs



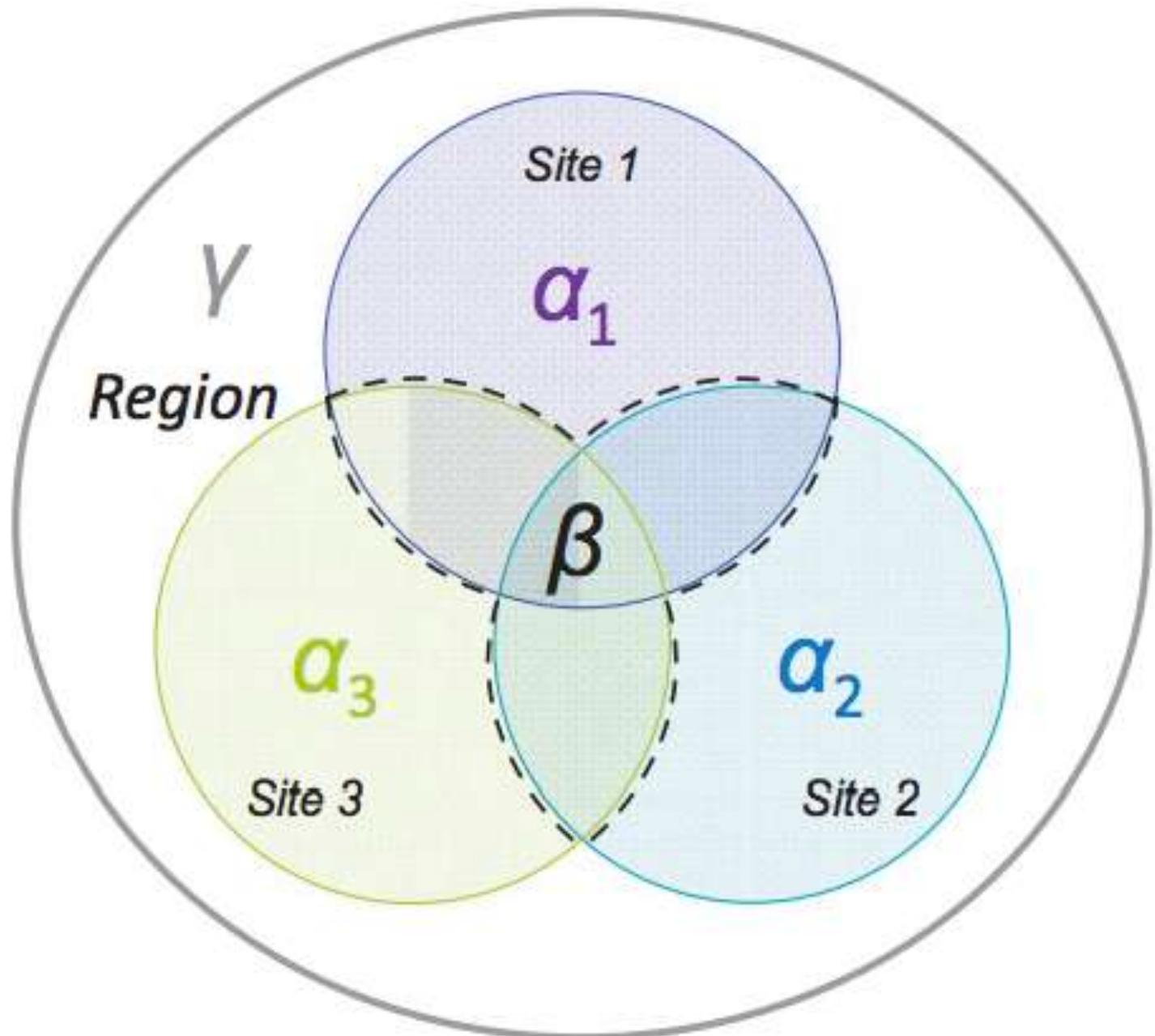
How do we compare between different samples?

Concept: Diversity measures

Measures of biodiversity

Zinger et al (2012)

“... measuring biodiversity consists of characterizing the **number, composition** and **variation** in taxonomic or functional units (**OTU**) over a wide range of biological organizations”



Measures of biodiversity

Zinger et al (2012)

Alpha diversity refers to the diversity within one location or sample. It is often measured as species richness (i.e. number of species), seldom as species evenness (extent of species dominance). Species richness is strongly sensitive to sampling effort, and requires standardized samples, or the use of estimators that corrects undersampling biases, such as Chao1 or ACE. Evenness is less affected by undersampling biases and is usually assessed with Simpson's or Pielou's indices or rank abundance curves (review in Magurran 2004).

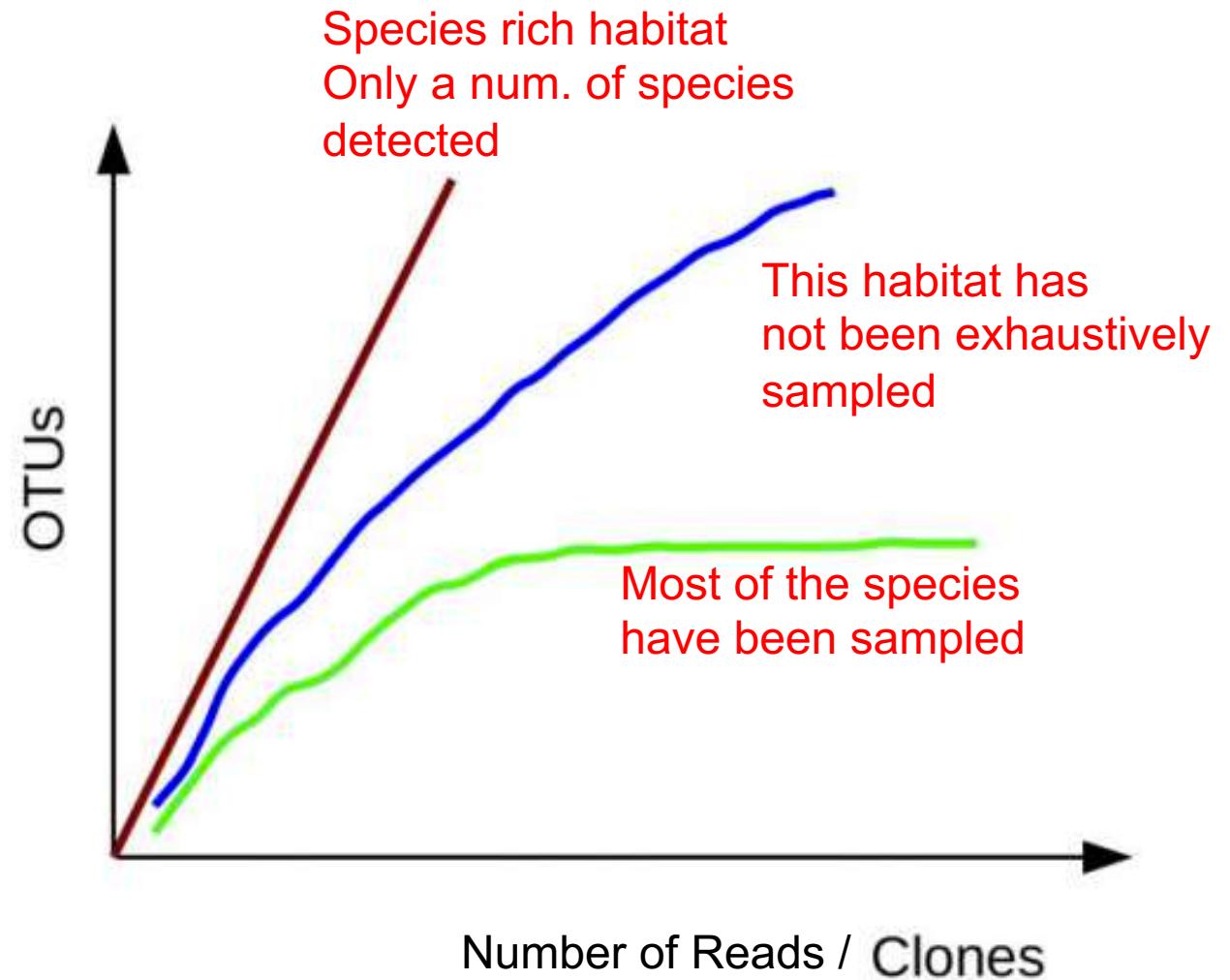
Beta diversity consists in determining the difference in diversity or community composition between two or more locations or samples (i) by considering species composition only, and use incidence data with associated metrics such as Jaccard or Sorensen similarity indices or (ii) by taking species relative abundances into account, and use Bray–Curtis or Morisita–Horn dissimilarity measures (Anderson *et al.* 2011). Using abundance data is, however, strongly discussed among microbiologists when dealing with rRNA gene data because of variations in gene copy number among strains (Acinas *et al.* 2004b; Zhu *et al.* 2005) as well as PCR artefacts.

Gamma diversity, or regional diversity, is similar to alpha diversity but applies for a larger area that encompasses the units under study.

Finally, the spatial scale of investigation can produce very different results and should be consistent in cross-study comparisons (Magurran 2004).

Species sampling and Rarefaction

Rarefaction allows the calculation of **species richness** for a given number of individual samples, based on the construction of so-called **rarefaction curves**. This curve is a plot of the number of species as a function of the number of samples



Alpha diversity

a measure of the diversity within a single sample

Types of alpha diversity

Total # of species = **richness**

How many OTUs?

Total # of genes = genetic richness

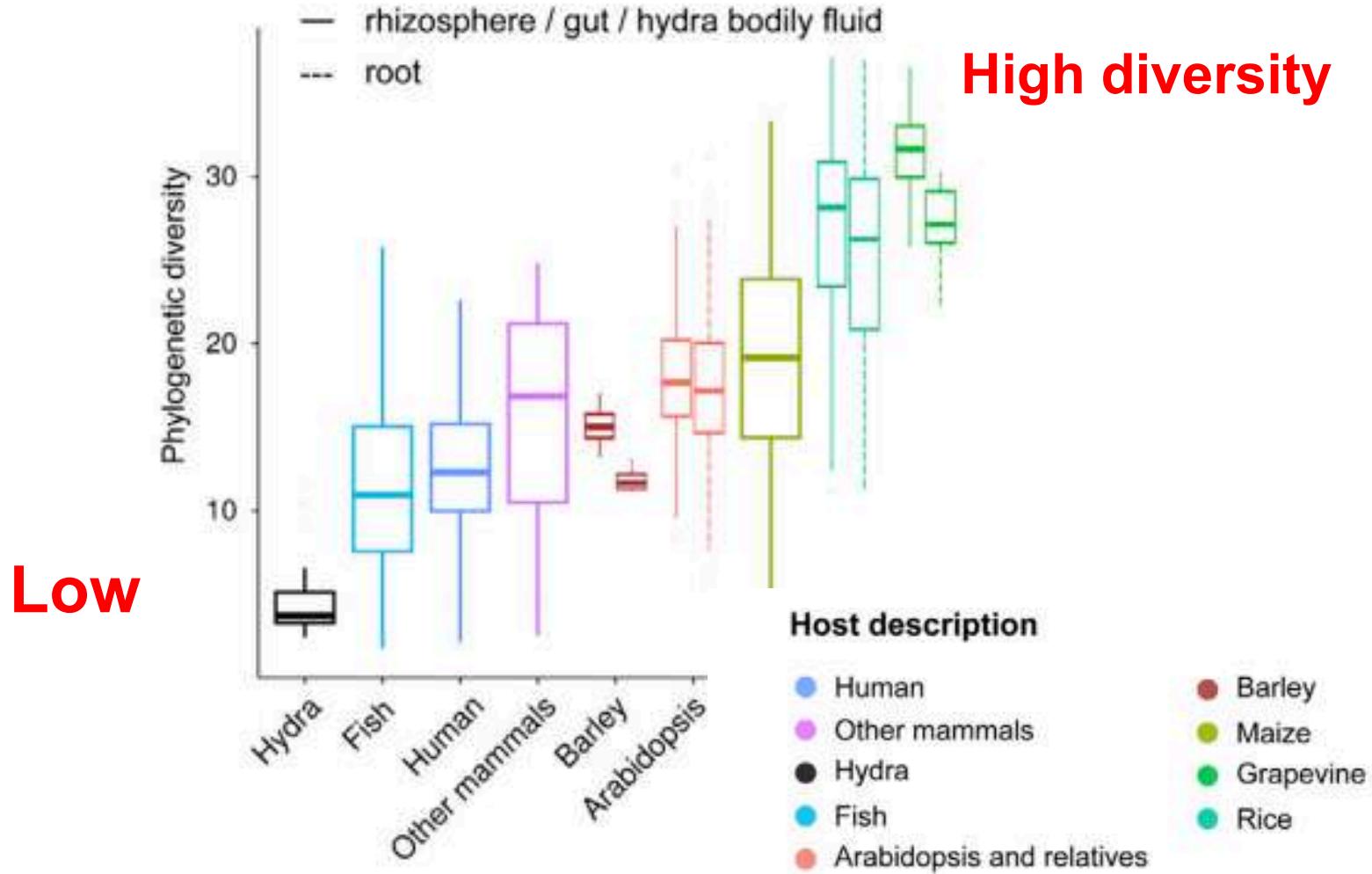
Phylogenetic diversity of genes = genetic PD

Evenness = What is the distribution of abundance in the community?

How many OTUs at high abundance and how many OTU at low abundance?

B

Alpha-diversity (phylogenetic diversity)



Beta diversity

a measure of **the similarity in diversity between samples**

Types of beta diversity

Species presence/absence

Shared phylogenetic diversity

Gene presence / absence

Shared phylogenetic diversity of genes

Frequently used as values for PCA of PCoA analysis

Beta diversity

A. Membership:

shared OTU occurrences across communities
1 = present, 0 = below detection

List of observed OTUs	Occurrences in community A	Occurrences in community B	Shared occurrences A & B
	OTU 1	0	X
OTU 2	0	1	
OTU 3	1	1	X
OTU 4	1	1	X
OTU 5	1	1	X

B. Composition:

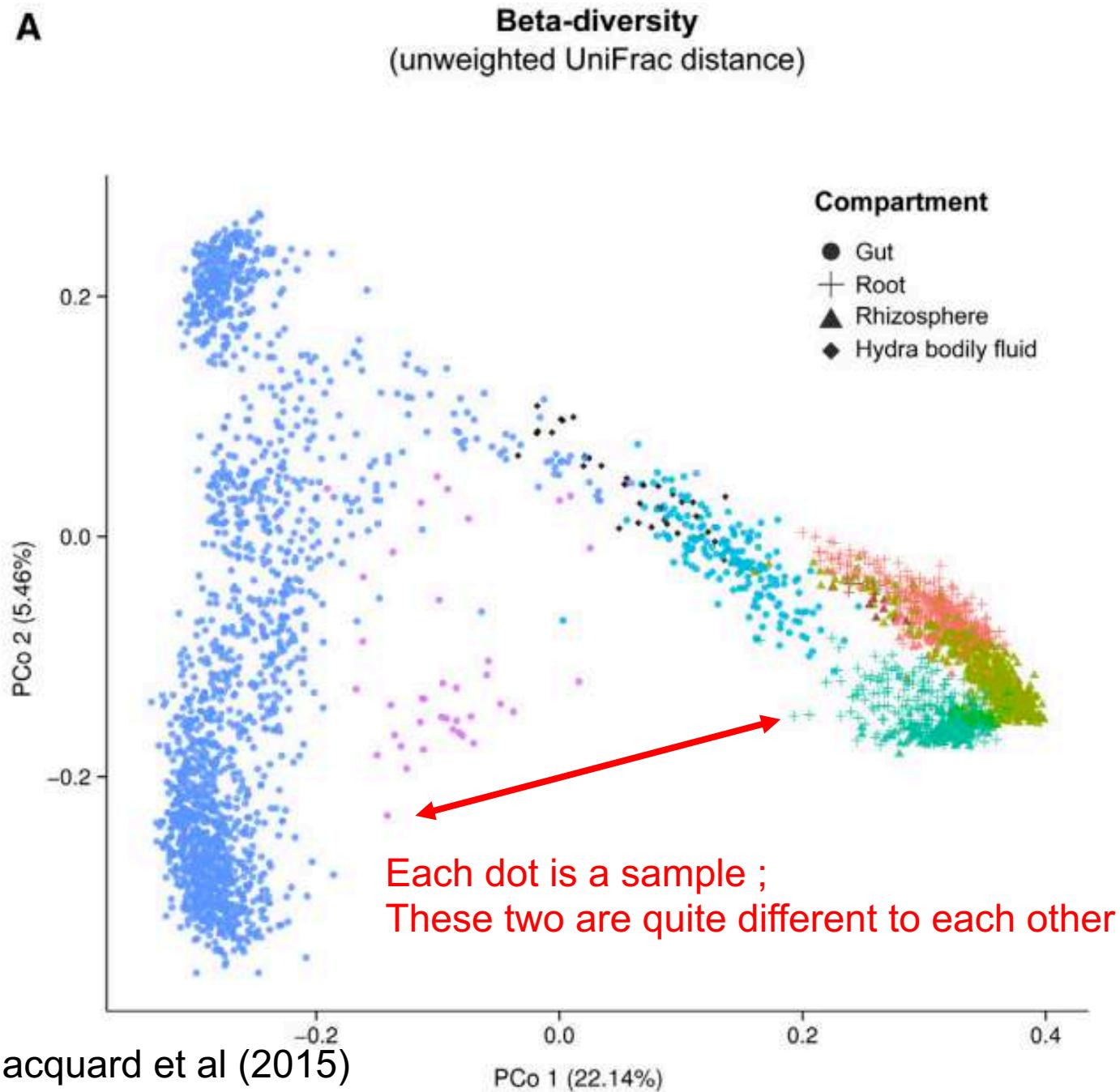
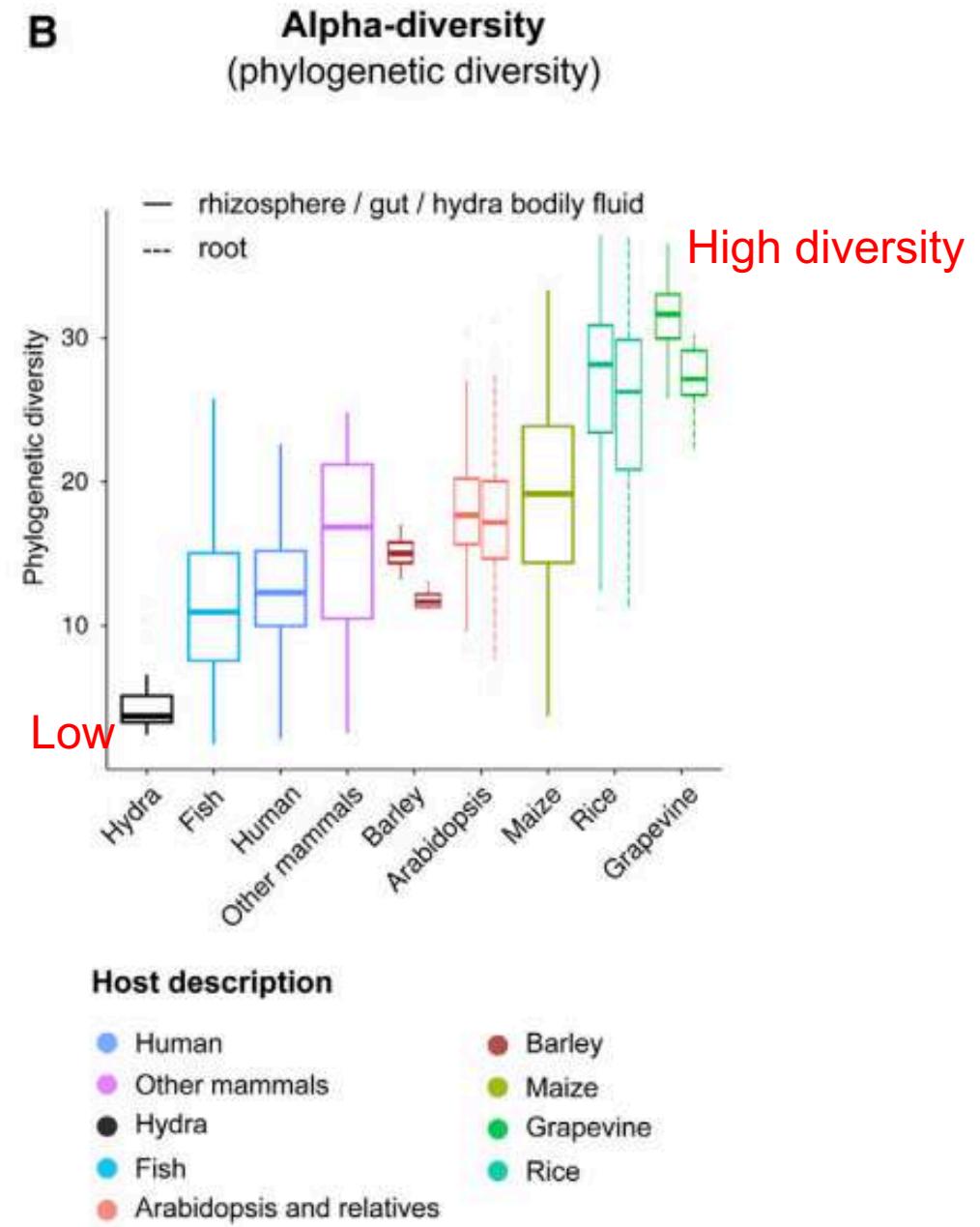
similar OTU abundances across communities

List of observed OTUs	Abundances community A	Abundances community B	Similar abundances A & B
	OTU 1	0.4	0
OTU 2	0	0.1	
OTU 3	0.1	0.1	X
OTU 4	0.2	0.5	
OTU 5	0.3	0.3	X

Phylogeny:

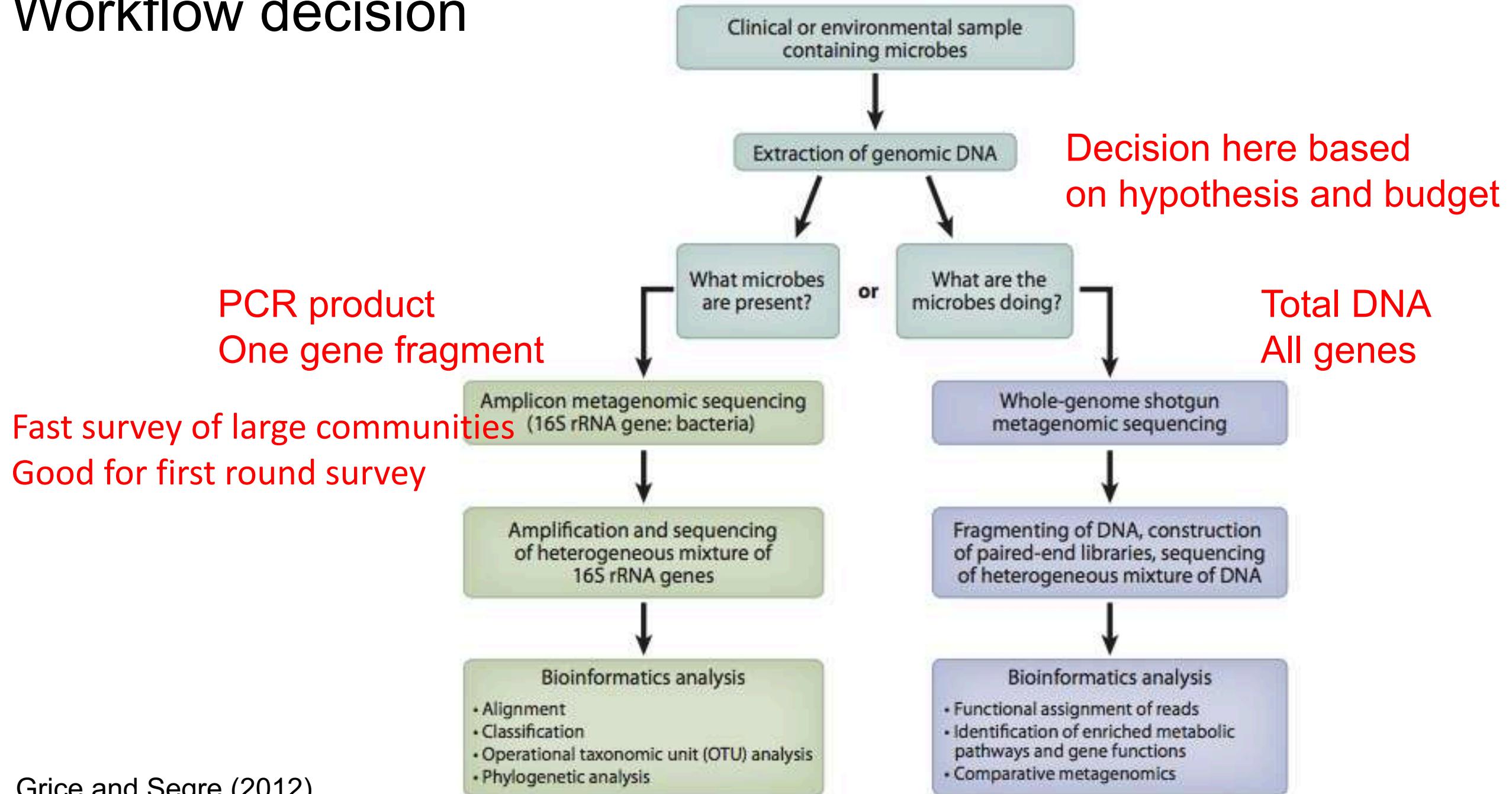
shared OTU lineages across communities

lineage	Abundances community A	Abundances community B	Similar abundances A & B
i.	OTU 1 0.4	0	X
ii.	OTU 2 0	0.1	
iii.	OTU 3 0.1	0.1	
	OTU 4 0	0.8	X
	OTU 5 0.5	0	

A**B**

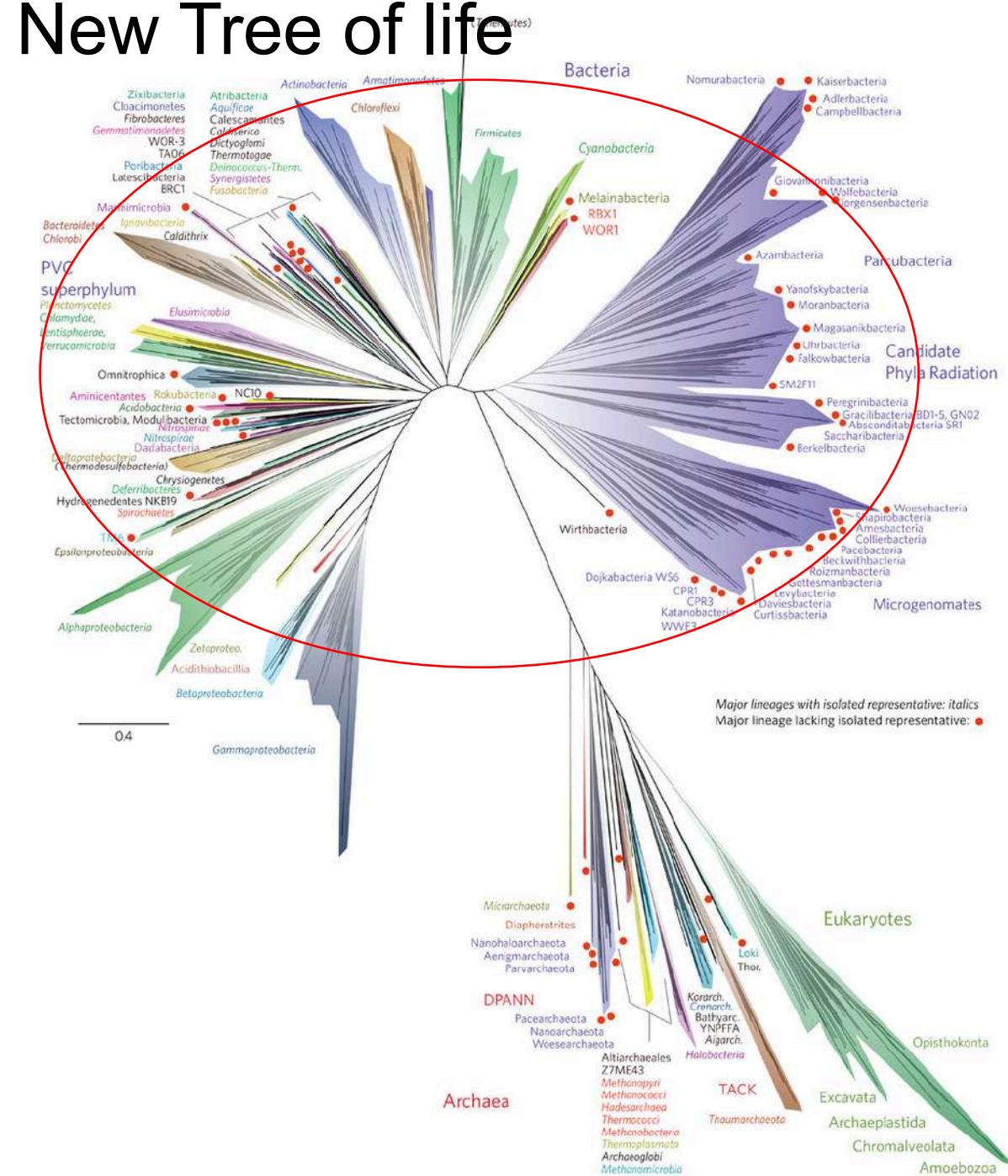
Amplicon sequencing or metagenomes?

Workflow decision



Amplicon sequencing

New Tree of life

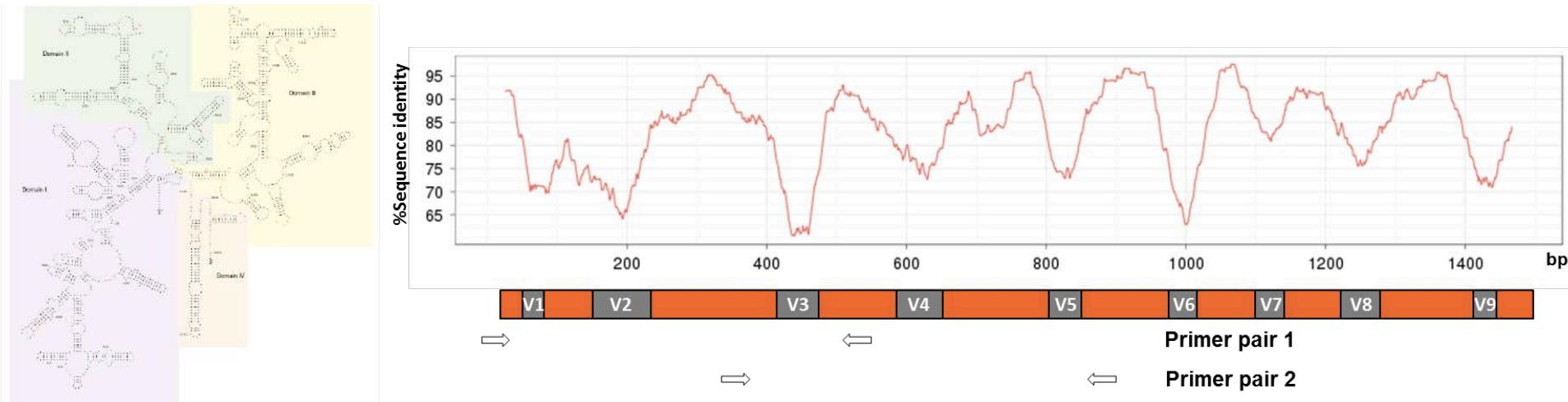


What do they have in common?

Hug et al (2016)

http://www.nytimes.com/2016/04/12/science/scientists-unveil-new-tree-of-life.html?_r=0

16S



- **Advantages:**

- Universal: Every bacterial and archaea species has this gene
- Conserved regions (for primer design)
- **Variable regions (to distinguish different species)**
- Great databases and alignments (for human related species)
- Mainly used for taxonomical classification

- **Problems:**

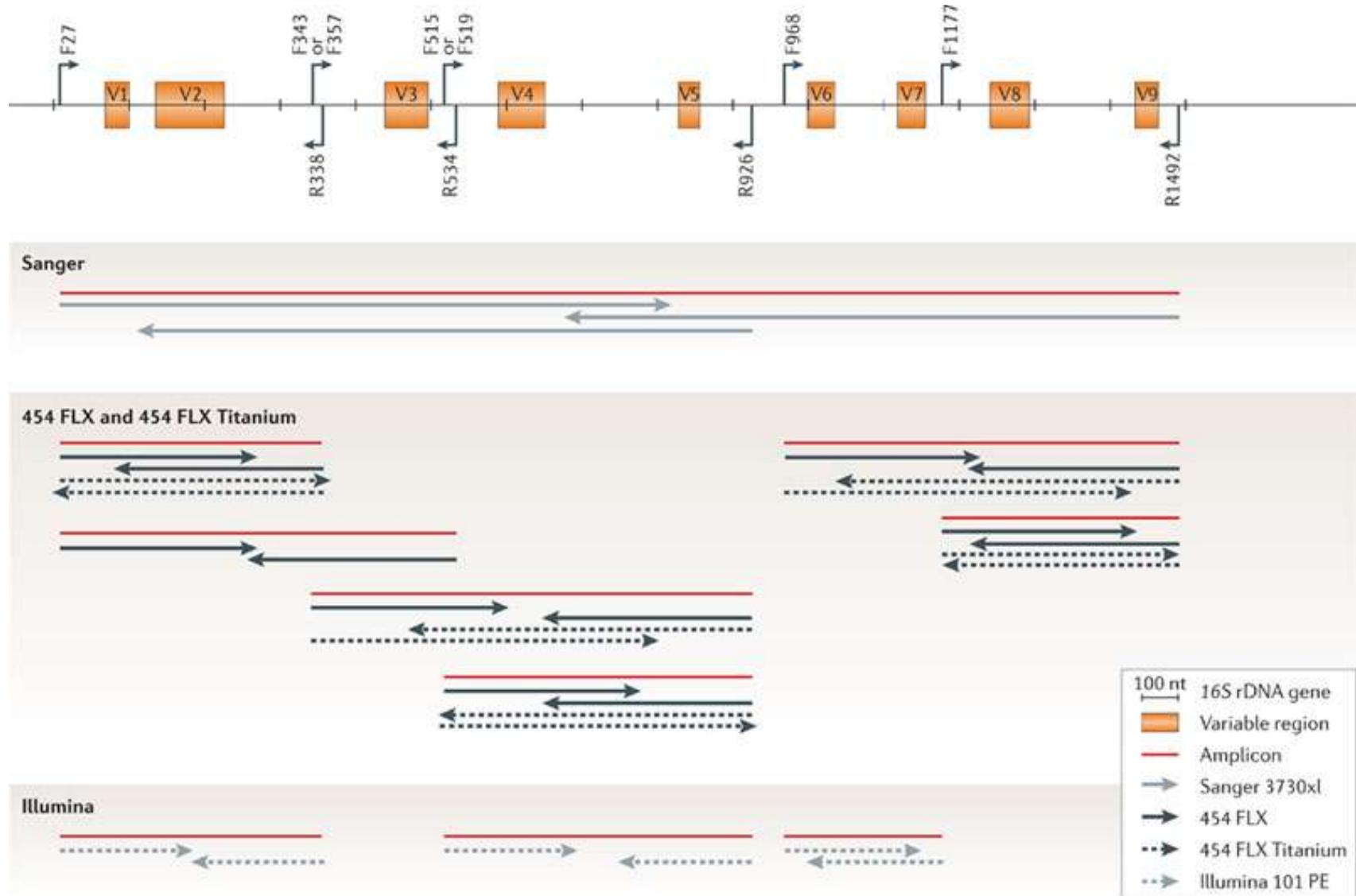
- Variable copy number in each species
- No universal (unbiased) primers
- (Not directly correlated with activity)
- (Lack of functional information)

Typical workflow

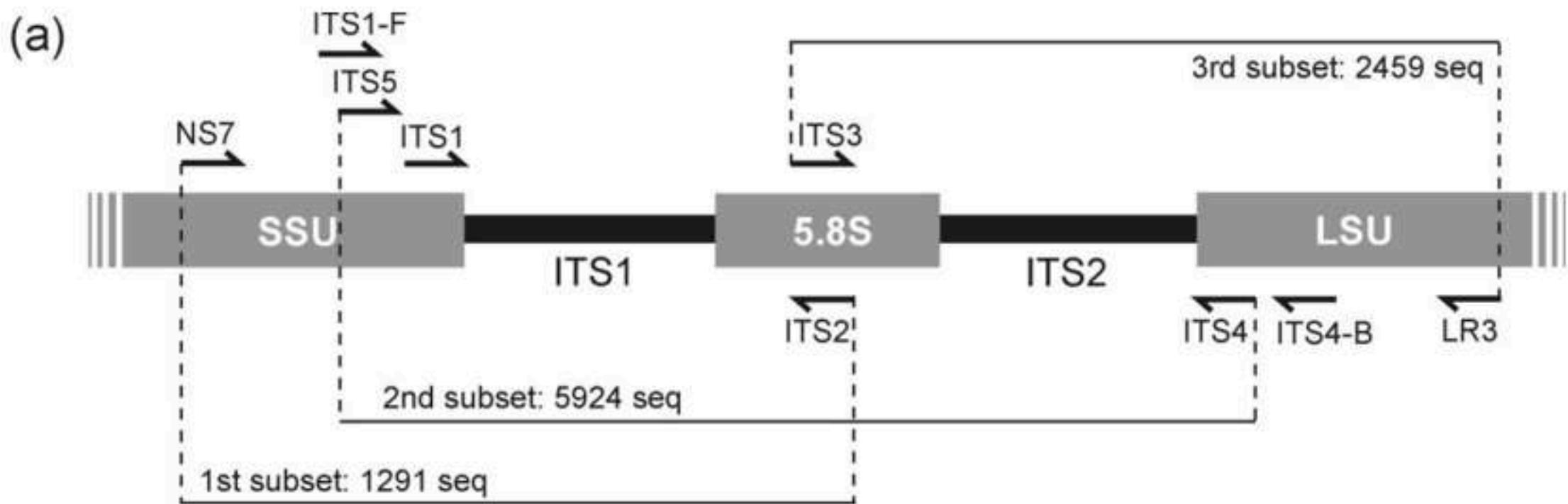


Which region to sequence?

16S amplified region



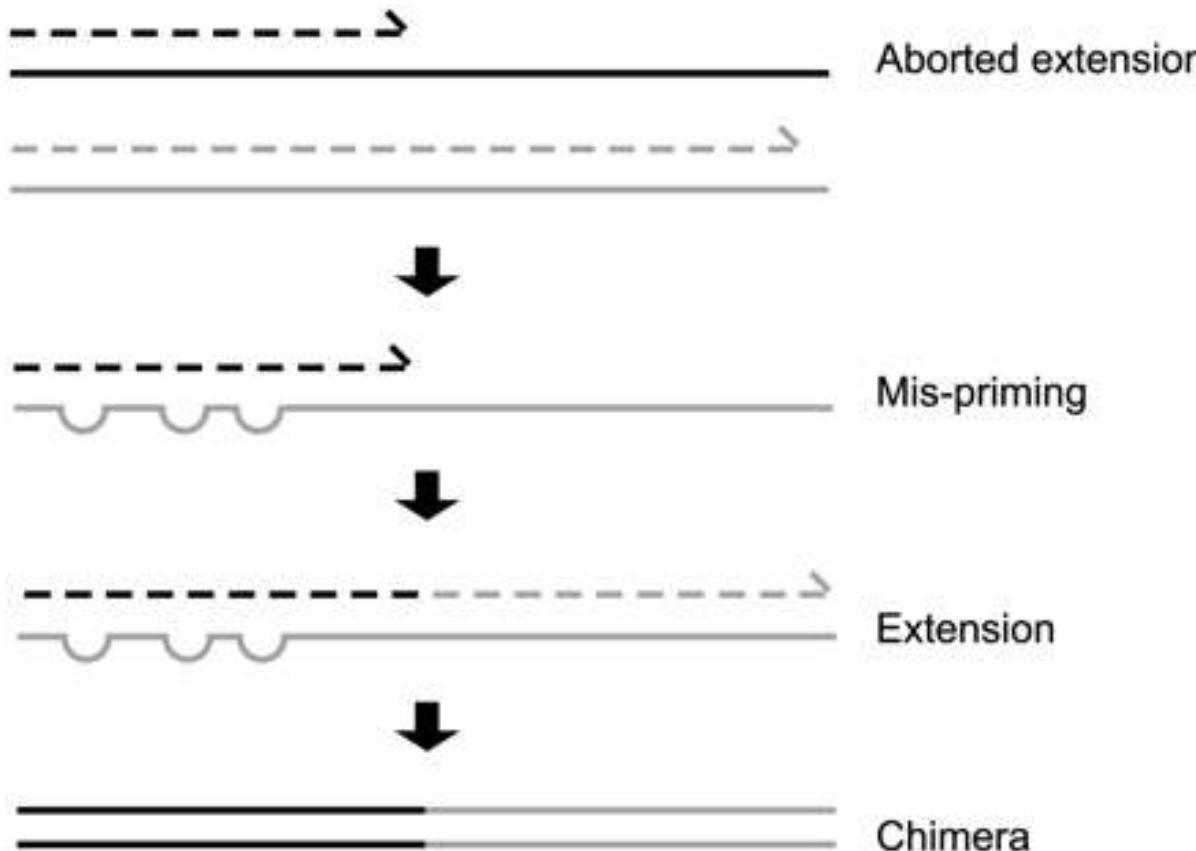
ITS for characterization of fungi species



Potential problem

- Lack of tools for processing ITS/Fungal microbiome data sets
- Amplification bias effects accuracy and replication
- Use of short reads prevents disambiguation of similar strains
- 16S or ITS may not differentiate between similar strains –
 - Clustering is done at 97%
 - Regions may be >99% similar
- Sequencing error inflates number of OTUs
- Chloroplast 16S sequences can get amplified in plant metagenomes

Chimeric 16S (Artificial sequences formed during PCR amplification)



“Chimeras were found to reproducibly form among independent amplifications and contributed to false perceptions of sample diversity and the false identification of novel taxa, **with less-abundant species exhibiting chimera rates exceeding 70%**”

Metagenomics

Advantage of metagenomics approach

Better classification with Increasing number of complete genomes

Focus on whole genome based phylogeny (whole genome phlyotyping)

- Advantages

No amplification bias like in 16S/ITS

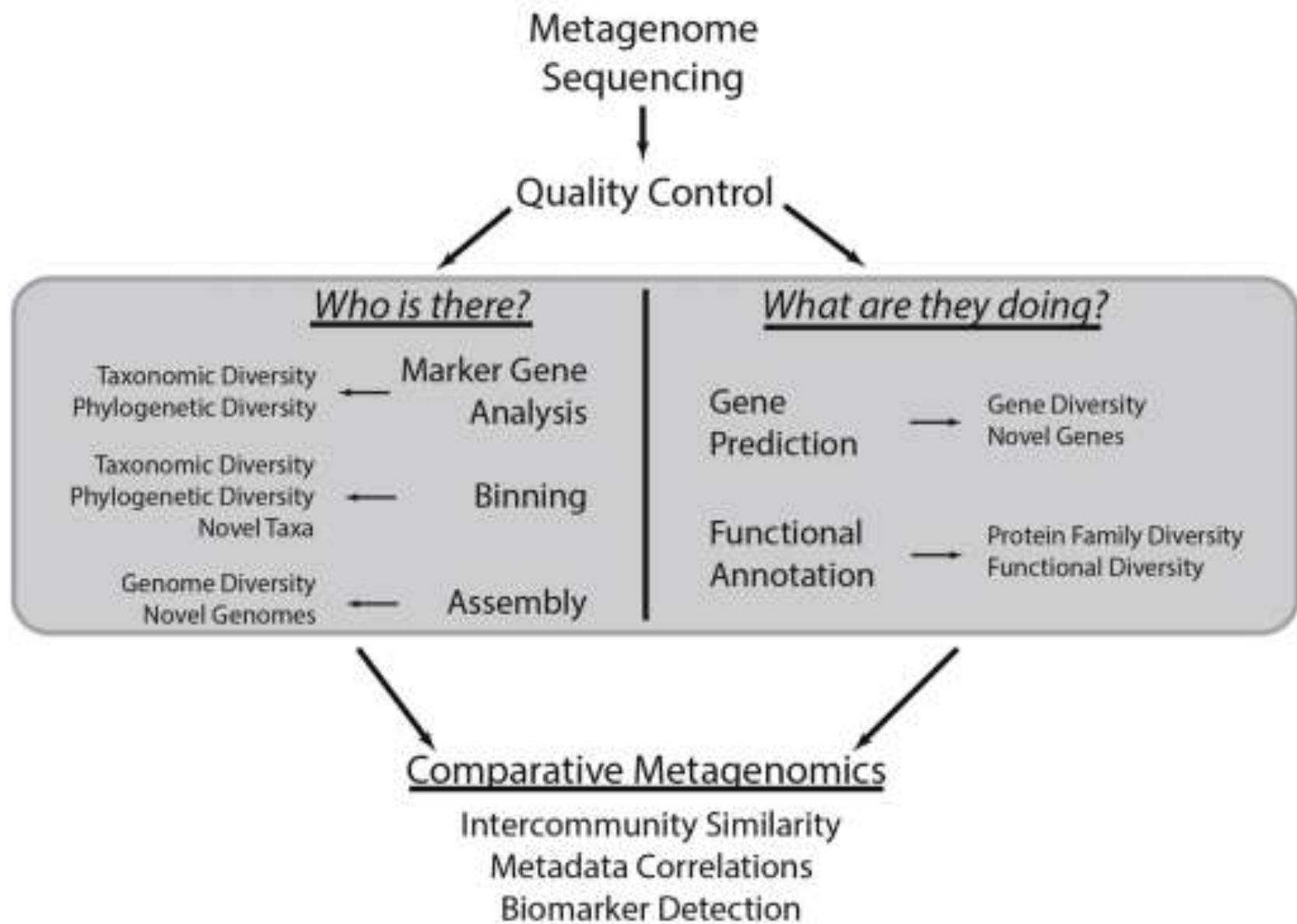
Issues

Poor sampling beyond eukaryotic diversity

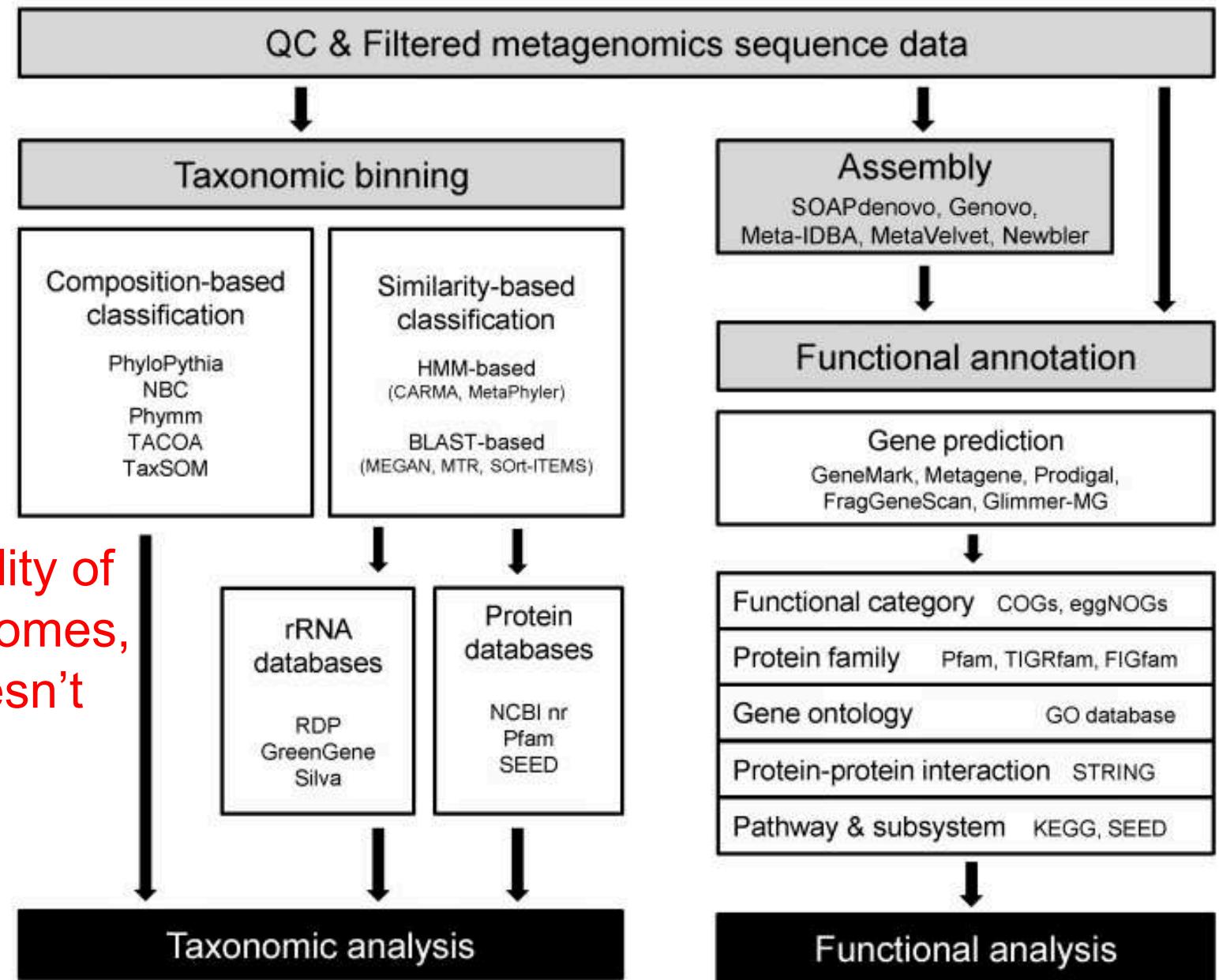
Assembly of metagenomes is **challenging** due to uneven coverage

Requires **high** depth of coverage

Overall workflow

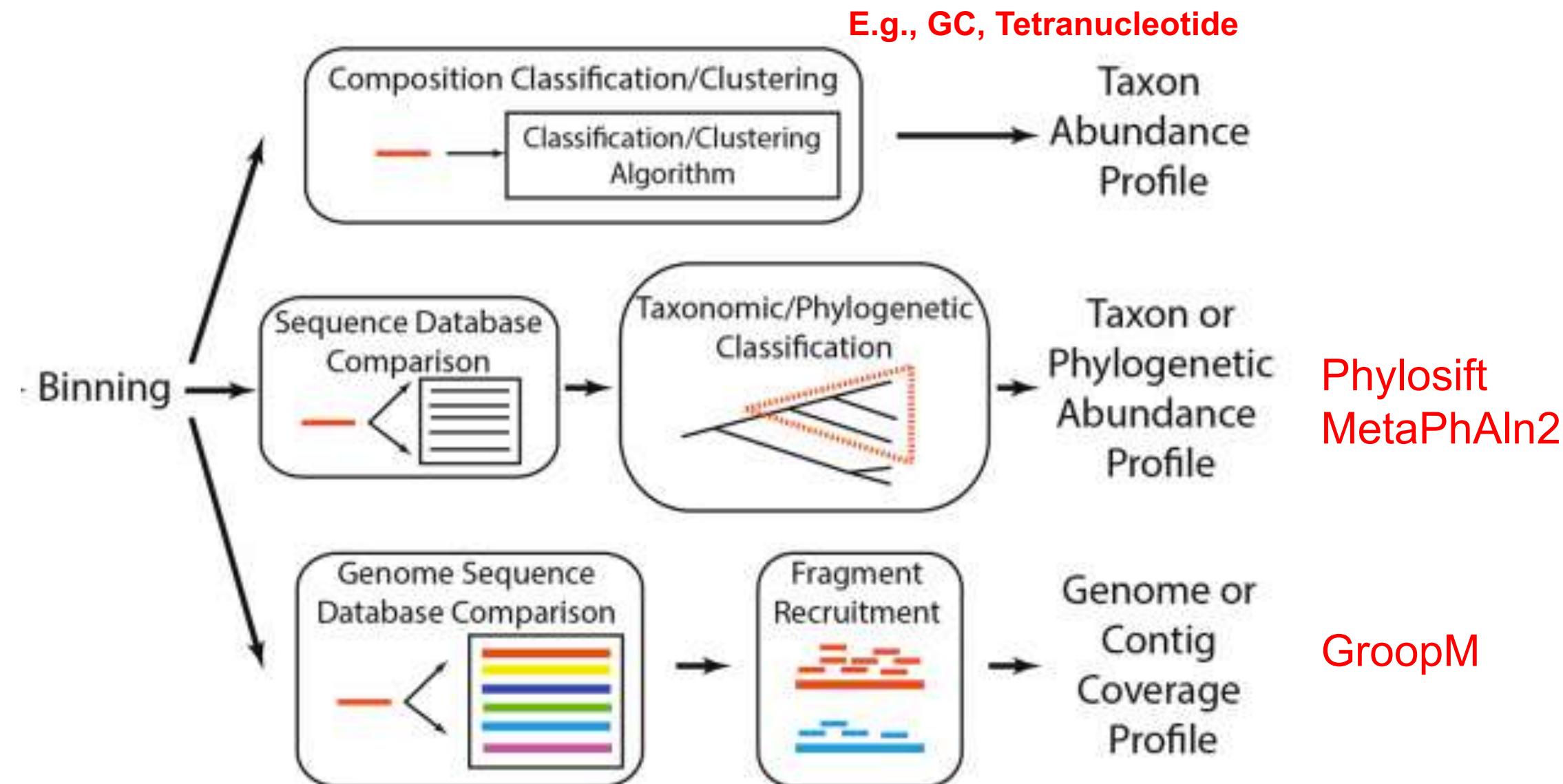


Overall workflow



With the increase availability of reference sequenced genomes, probably one day one doesn't require assembly of metagenomes

Binning methods



Binning methods: A combination of

Classification based on **sequence composition**:

Advantage : all reads can be categorised into bins

Disadvantage: no taxonomy / function of the bins.

Classification based on **sequence similarity (of known genes)**

Advantage: One can determine taxonomy and function of reads.

Disadvantage: reads with similarity can not be classified .

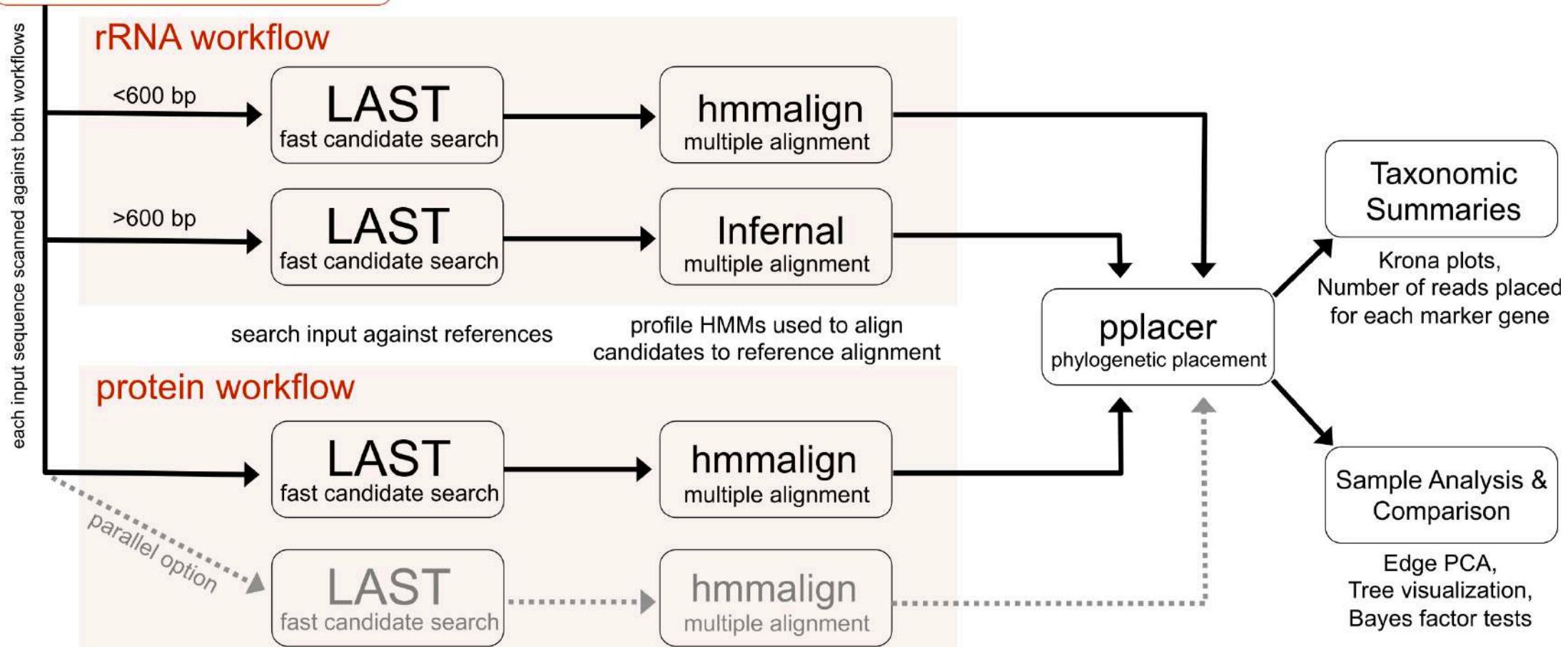
PhyloSift

mining the global metagenome

<https://phylosift.wordpress.com/>

- Uses a database of 37 universal proteins & rRNA genes.
- Designed to classify using phylogenies

Input Sequences



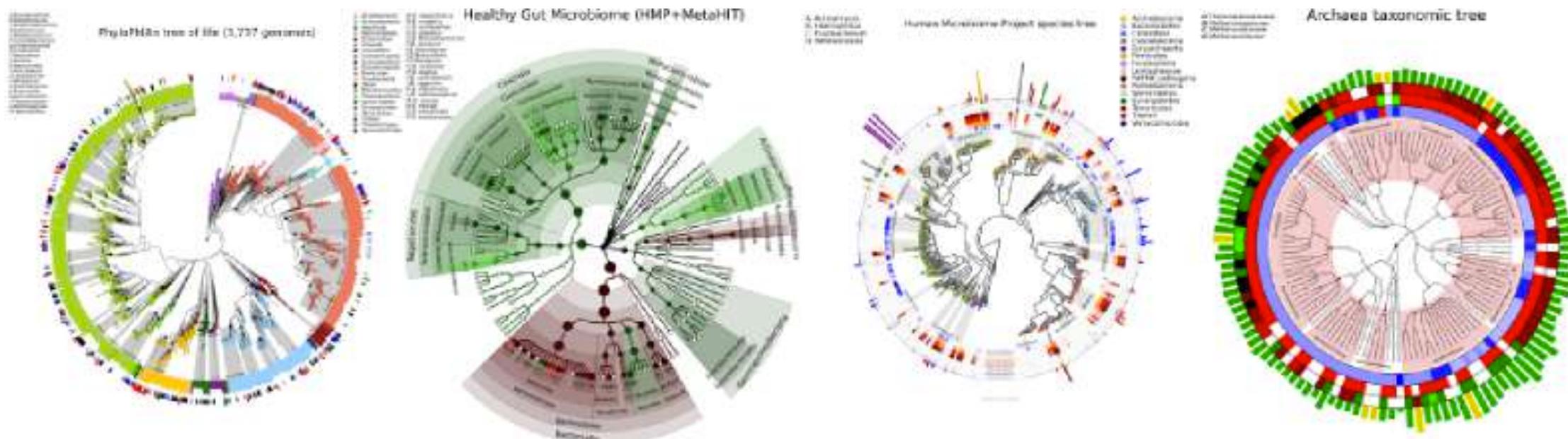
MetaPhAln2 – enhanced metagenomic taxonomic profiling

relies on ~1M unique clade-specific marker genes identified from ~17,000 reference genomes (~13,500 bacterial and archaeal, ~3,500 viral, and ~110 eukaryotic), allowing:

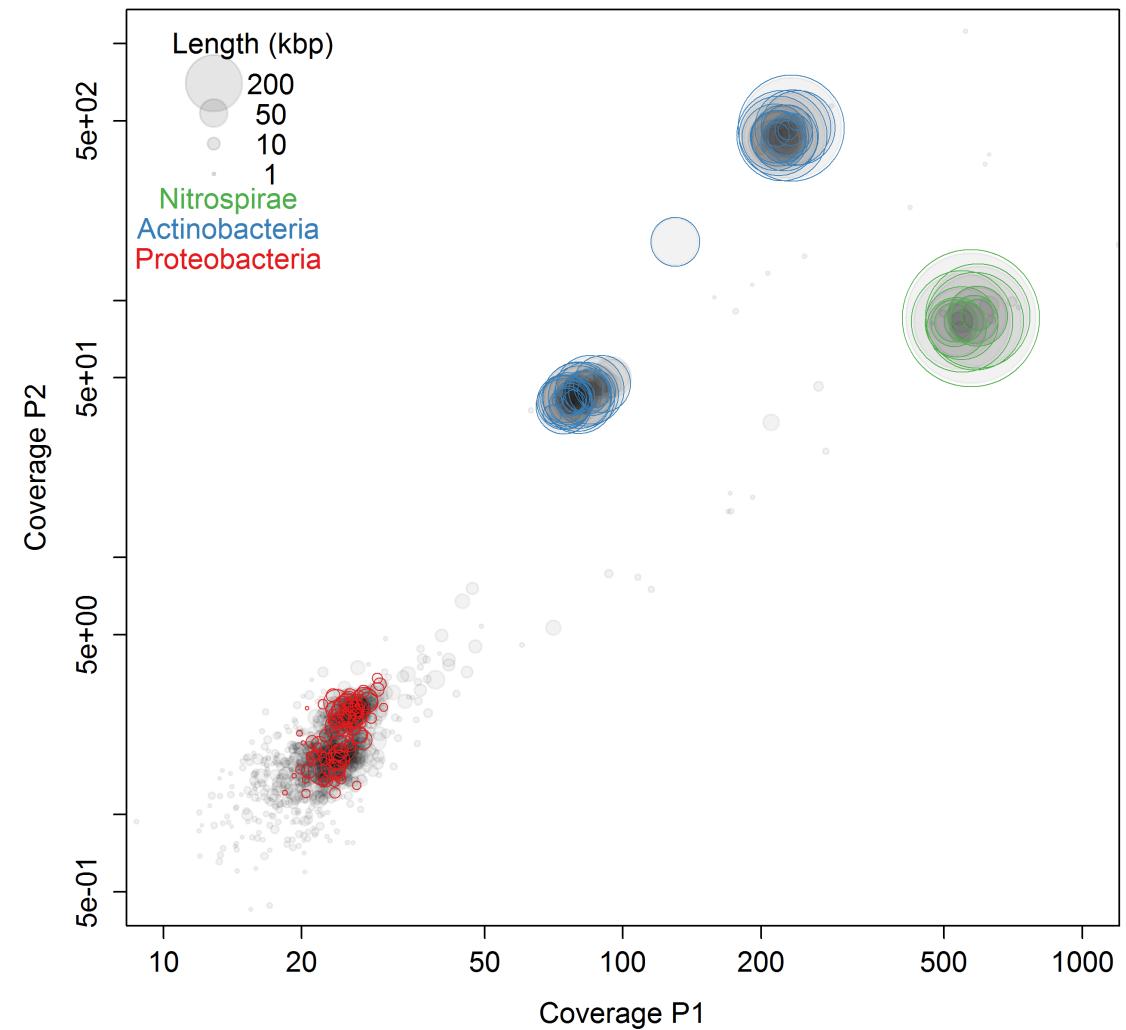
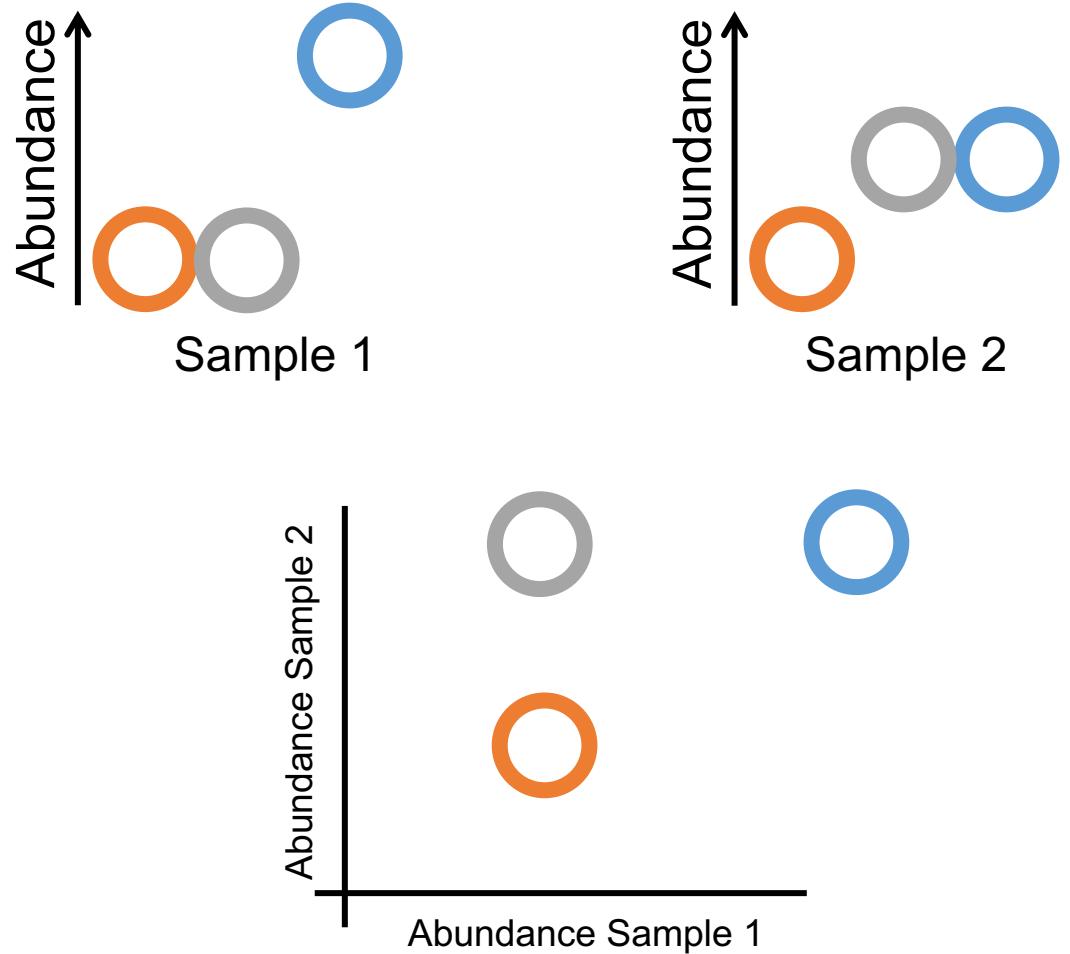
Species level resolution

Good visualisation with **GraphAln**

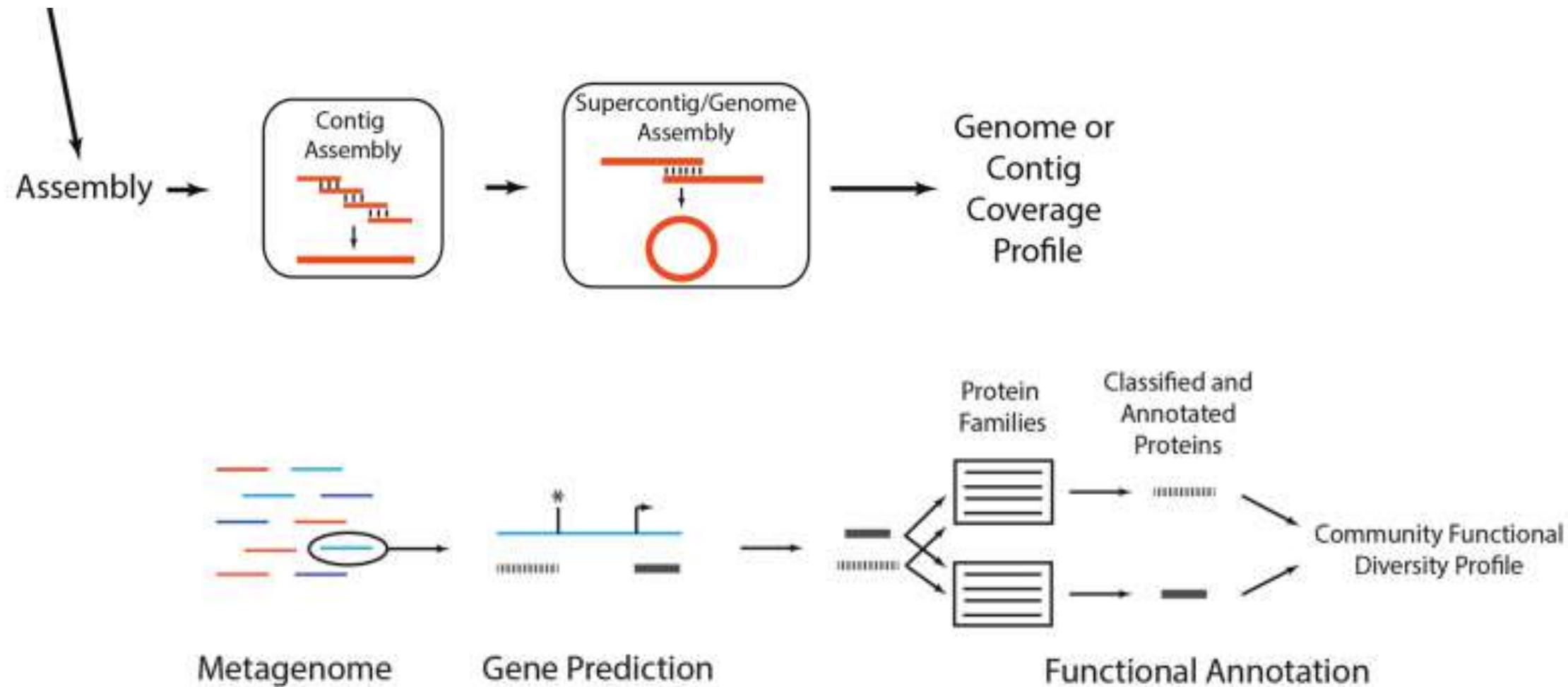
(So it's useful with known ecosystems)



Example of binning based on differential coverage



Actual assembly



Algorithm advancements lead to recovery of genomes

nature
microbiology

ARTICLES

DOI: 10.1038/s41564-017-0012-7

OPEN

Corrected: Author correction

Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life

Donovan H. Parks , Christian Rinke , Maria Chuvochina, Pierre-Alain Chaumeil, Ben J. Woodcroft, Paul N. Evans, Philip Hugenholtz * and Gene W. Tyson*

Biases

Extraction protocol matters



Soil Biology & Biochemistry 36 (2004) 1607–1614

Soil Biology &
Biochemistry

www.elsevier.com/locate/soilbio

Impact of DNA extraction method on bacterial community composition measured by denaturing gradient gel electrophoresis

Julia R. de Liphay^{a,b}, Christiane Enzinger^{b,1}, Kaare Johnsen^{a,2}, Jens Aamand^a,
Søren J. Sørensen^{b,*}

^aDepartment of Geochemistry, Geological Survey of Denmark and Greenland, Øster Voldgade 10, DK-1350 Copenhagen K, Denmark

^bDepartment of Microbiology, University of Copenhagen, Sølygade 8/H, DK-1307 Copenhagen K, Denmark

Received 1 September 2003; received in revised form 6 March 2004; accepted 15 March 2004

Abstract

The impact of DNA extraction protocol on soil DNA yield and bacterial community composition was evaluated. Three different procedures to physically disrupt cells were compared: sonication, grinding-freezing-thawing, and bead beating. The three protocols were applied to three different topsoils. For all soils, we found that each DNA extraction method resulted in unique community patterns as measured by denaturing gradient gel electrophoresis. This indicates the importance of the DNA extraction protocol on data for evaluating soil bacterial diversity. Consistently, the bead-beating procedure gave rise to the highest number of DNA bands, indicating the highest number of bacterial species. Supplementing the bead-beating procedure with additional cell-rupture steps generally did not change the bacterial community profile. The same consistency was not observed when evaluating the efficiency of the different methods on soil DNA yield. This parameter depended on soil type. The DNA size was of highest molecular weight with the sonication and grinding-freezing-thawing procedures (approx. 20 kb). In contrast, the inclusion of bead beating resulted in more sheared DNA (approx. 6–20 kb), and the longer the bead-beating time, the higher the fraction of low-molecular weight DNA. Clearly, the choice of DNA extraction protocol depends on soil type. We found, however, that for the analysis of indigenous soil bacterial communities the bead-beating procedure was appropriate because it is fast, reproducible, and gives very pure DNA of relatively high molecular weight. And very importantly, with this protocol the highest soil bacterial diversity was obtained. We believe that the choice of DNA extraction protocol will influence not only the determined phylogenetic diversity of indigenous microbial communities, but also the obtained functional diversity. This means that the detected presence of a functional gene—and thus the indication of various activities—may depend on the nature of the applied DNA extraction procedure.

“we found that each DNA extraction method resulted in unique community patterns”

Wesolowska-Andersen et al. *Microbiome* 2014, 2:19
<http://www.microbiomejournal.com/content/2/1/19>



Microbiome

Open Access

RESEARCH

Choice of bacterial DNA extraction method from fecal material influences community structure as evaluated by metagenomic analysis

Agata Wesolowska-Andersen¹, Martin Iain Bah², Vera Carvalho², Karsten Kristiansen³, Thomas Sicheritz-Pontén¹, Ramneek Gupta^{1*} and Tine Rask Licht^{2*}

Abstract

Background: In recent years, studies on the human intestinal microbiota have attracted tremendous attention. Application of next generation sequencing for mapping of bacterial phylogeny and function has opened new doors to this field of research. However, little attention has been given to the effects of choice of methodology on the output resulting from such studies.

Results: In this study we conducted a systematic comparison of the DNA extraction methods used by the two major collaborative efforts: The European MetaHIT and the American Human Microbiome Project (HMP). Additionally, effects of homogenizing the samples before extraction were addressed. We observed significant differences in distribution of bacterial taxa depending on the method. While eukaryotic DNA was most efficiently extracted by the MetaHIT protocol, DNA from bacteria within the Bacteroidetes phylum was most efficiently extracted by the HMP protocol.

Conclusions: Whereas it is comforting that the inter-individual variation clearly exceeded the variation resulting from choice of extraction method, our data highlight the challenge of comparing data across studies applying different methodologies.

“We observed significant differences in distribution of bacterial taxa depending on the method.”

Alpha diversity is always overestimated

Table 1. Effect of quality filtering and clustering on diversity estimates (OTU number), error rate and data loss of pyrotags amplified from two regions of *E. coli* MG1655 16S rRNA genes.

Read filtering	Number of OTUs at percentage identity thresholds						% errorless reads	% reads used
	100	99	98	97	95	90		
5' forward (V1 and V2)								
Theoretical number	5	4	3	1	1	1		
No quality filtering	643	95	31	16	5	3	68.7	77.9
Reads with N's removed	600	85	29	14	4	3	69.8	76.7
Quality score-based filtering (% per-base error probability)								
3	638	92	31	13	3	3	68.9	77.7
2	632	90	30	14	3	3	69.0	77.6
1	609	79	24	9	3	3	69.1	77.3
0.5	562	66	15	7	3	3	70.7	75.3
0.2	469	30	6	3	3	3	73.2	70.8
0.1	372	26	5	3	3	3	77.8	57.8
3' reverse (V8)								
Theoretical number	1	1	1	1	1	1		
No quality filtering	385	43	13	7	5	4	84.6	94.4
Reads with N's removed	361	40	12	6	4	3	85.3	93.6
Quality score-based filtering (% per-base error probability)								
3	378	40	12	7	5	4	84.8	94.2
2	368	32	10	6	5	4	85.1	93.8
1	342	25	9	6	5	4	85.3	93.3
0.5	310	20	8	6	5	4	87.5	89.5
0.2	236	7	2	2	2	2	89.6	82.1
0.1	196	4	2	2	2	2	90.7	70.6

Diversity estimates should be considered relative to the theoretical number of OTUs from *E. coli*.

Kunin et al (2010)

Reagent and laboratory contamination

RESEARCH ARTICLE

Open Access

Reagent and laboratory contamination can critically impact sequence-based microbiome analyses

Susannah J Salter^{1*}, Michael J Cox², Elena M Turek², Szymon T Calus³, William O Cookson², Miriam F Moffatt², Paul Turner^{4,5}, Julian Parkhill¹, Nicholas J Loman³ and Alan W Walker^{1,6*}

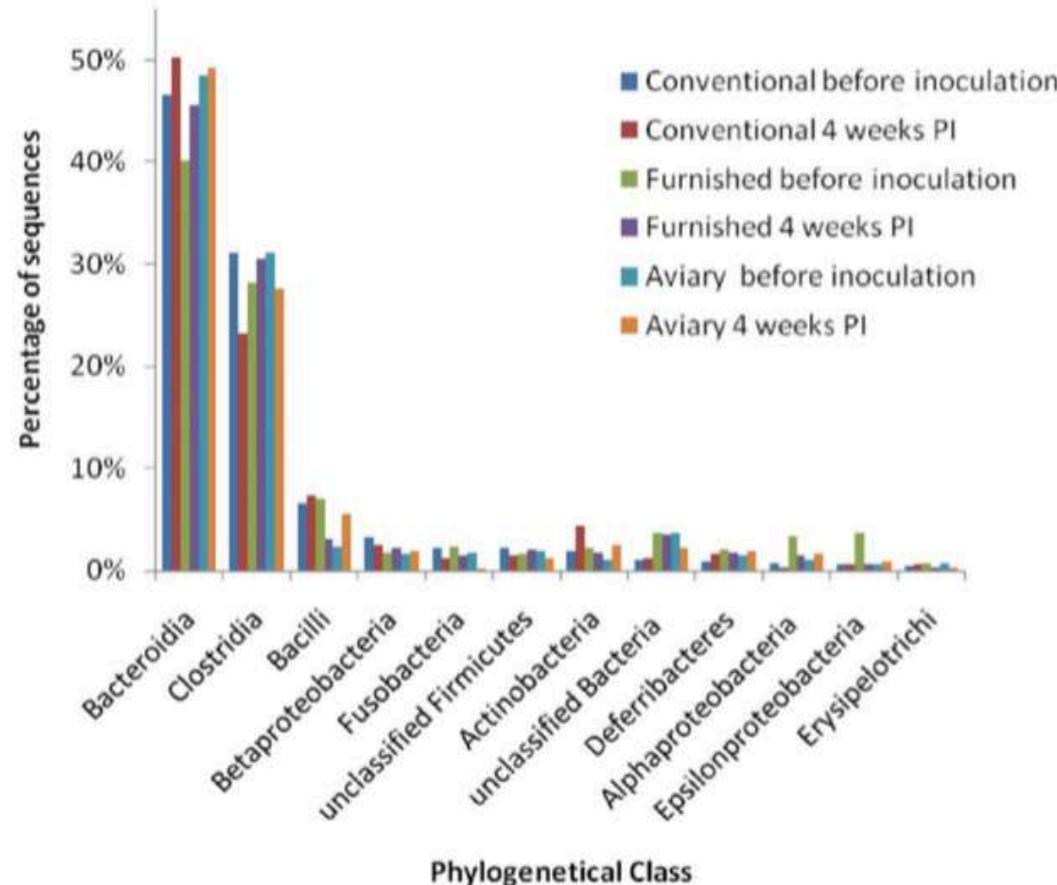
RESEARCH HIGHLIGHT

Tracking down the sources of experimental contamination in microbiome studies

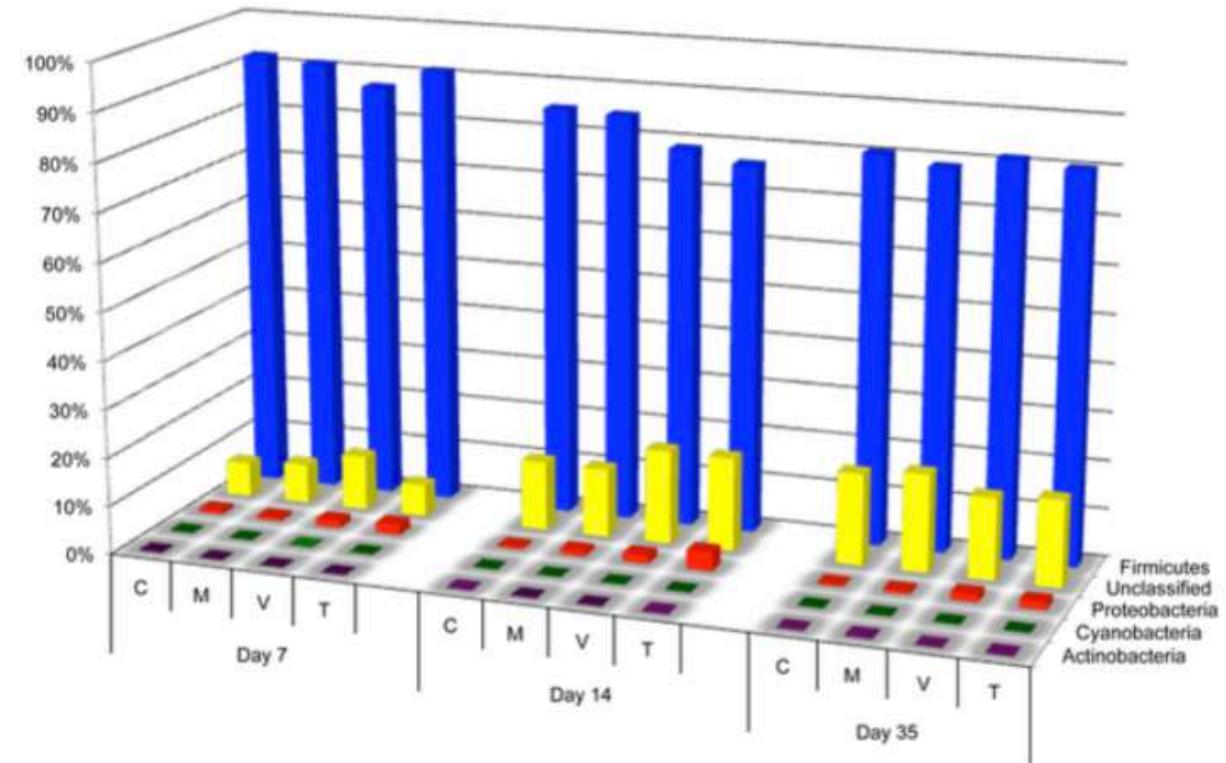
Sophie Weiss¹, Amnon Amir², Embriette R Hyde², Jessica L Metcalf², Se Jin Song² and Rob Knight^{2,3,4*}

2 papers with different results at the same year

Bacteroidetes >>> rest

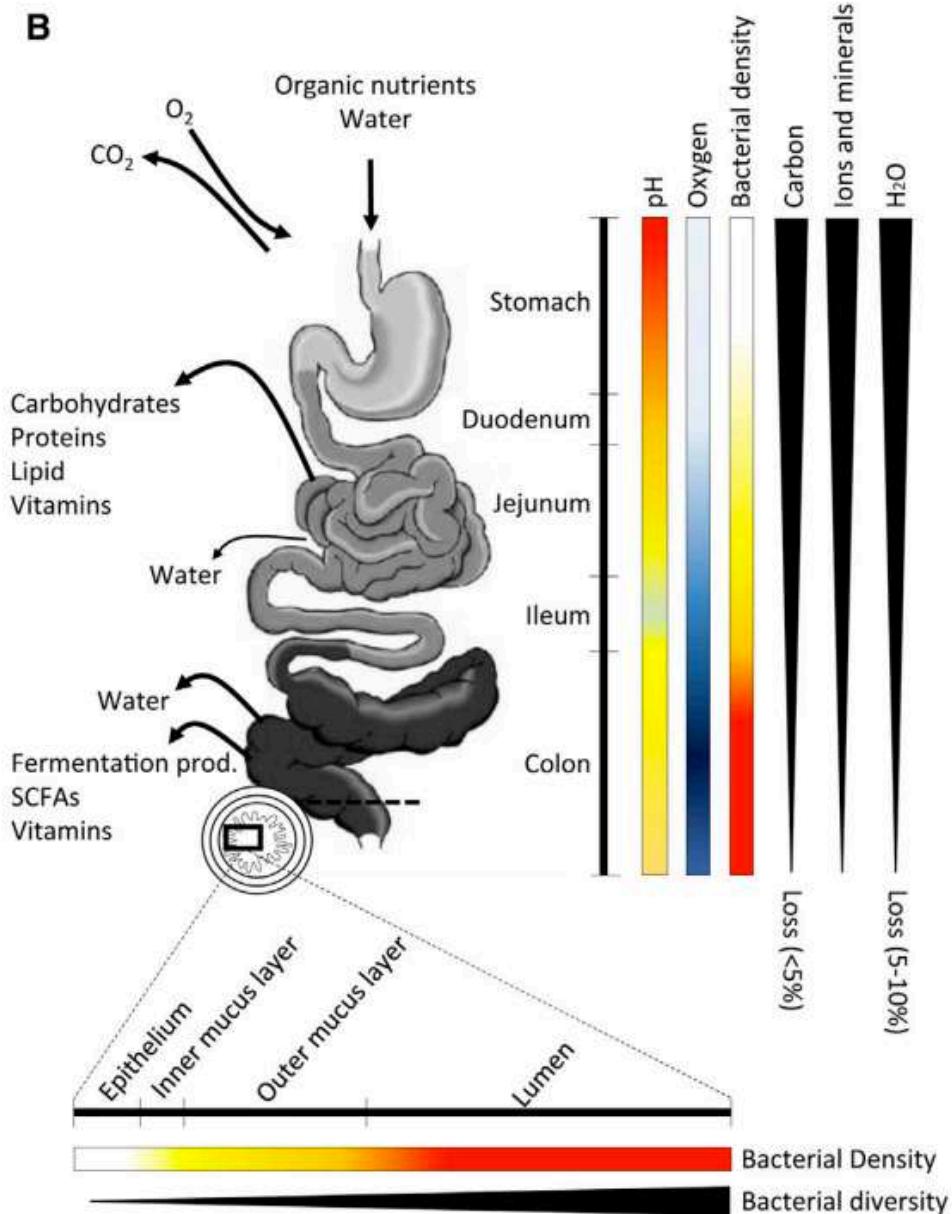
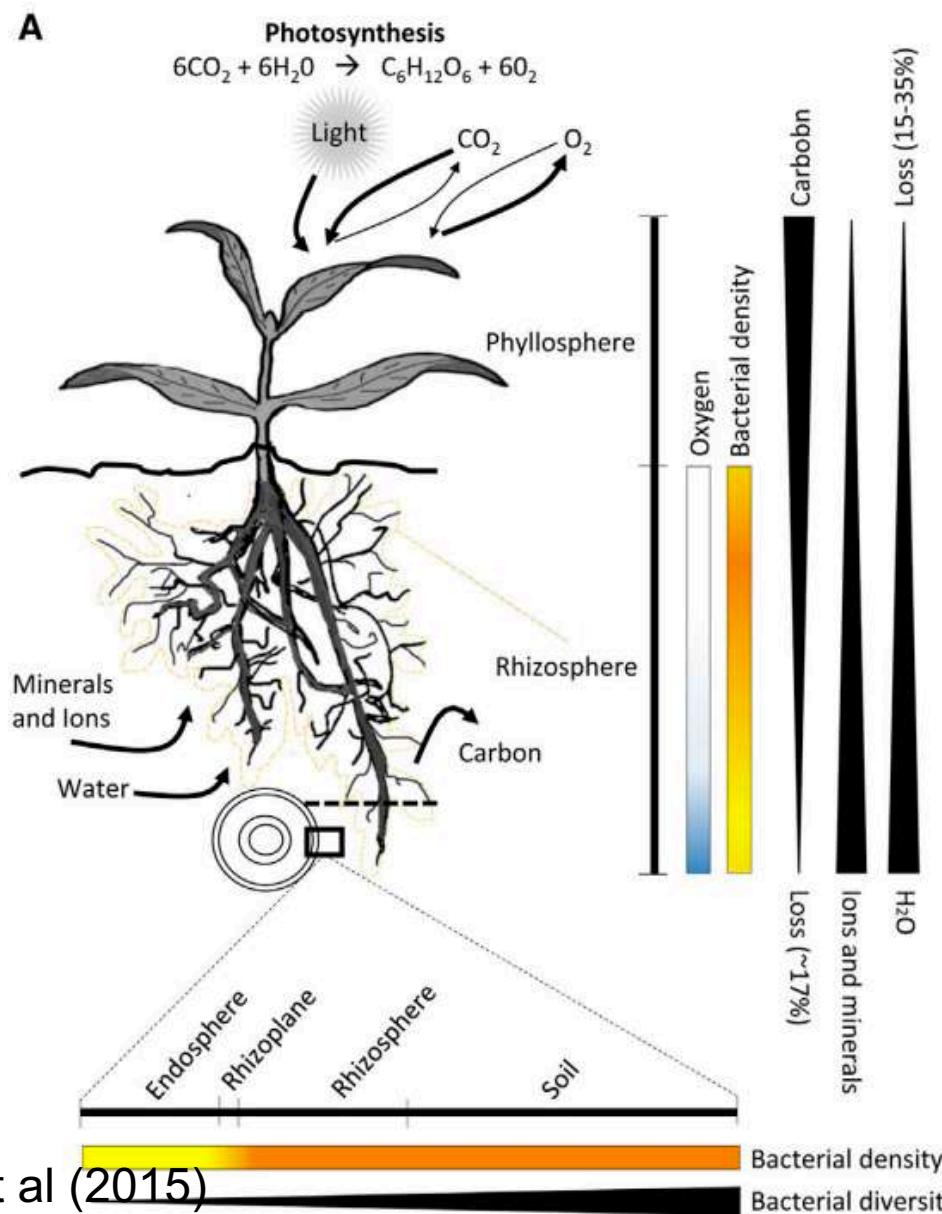


firmicutes >>> rest > bacteroidetes

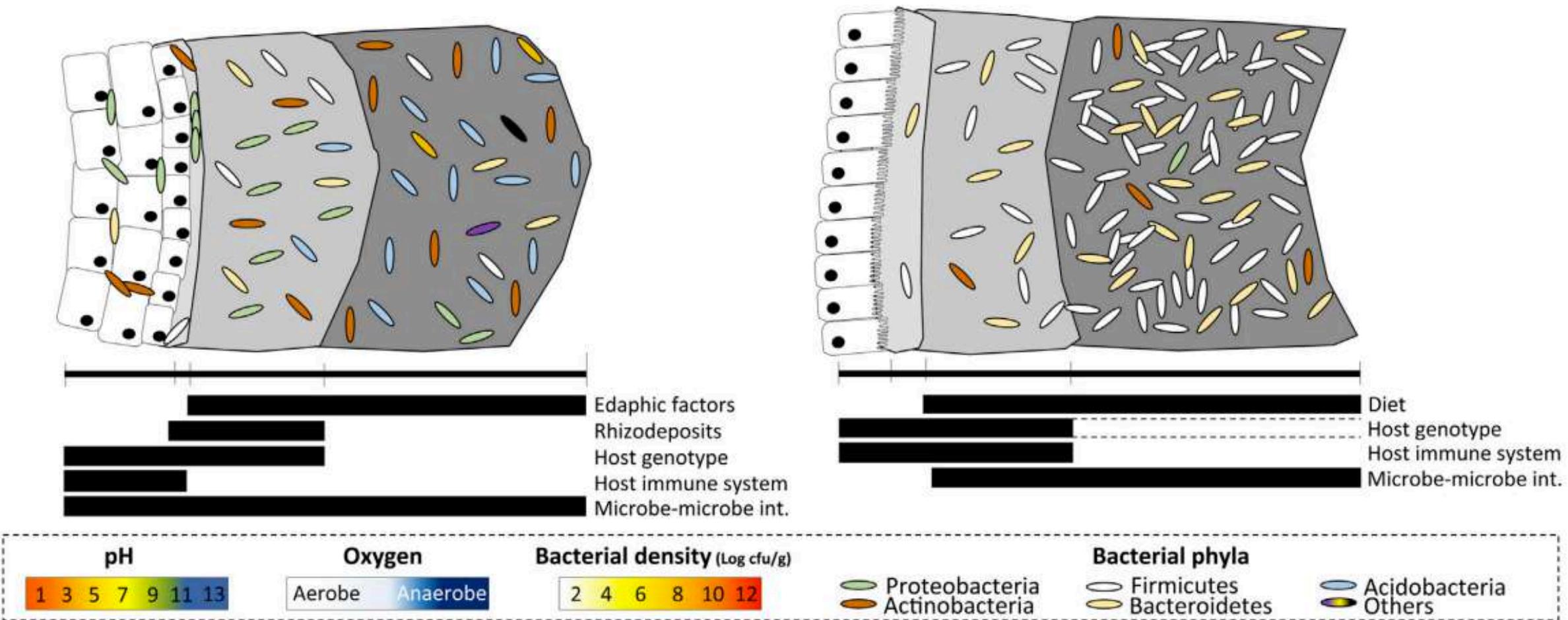


Case studies

Two most common systems



Two most common systems



Two most common systems

Table 1. Percentage of Shotgun Metagenome Reads Assigned to Each Kingdom of Life across Metagenome Studies

	Cucumber ^a	Wheat ^a	Soybean ^b	Wheat ^c	Oat ^c	Pea ^c	Barley ^d	Gut ^e
Bacteria	99.36	99.45	96	88.5	77.3	73.7	94.04	99.1
Archaea	0.02	0.02	<1	<0.5	<0.5	<0.5	0.054	
Eukaryotes	0.54	0.48	3	3.3	16.6	20.7	5.90	<0.1

^aOfek-Lalzar et al. (2014) (metagenomics of rhizoplane samples).

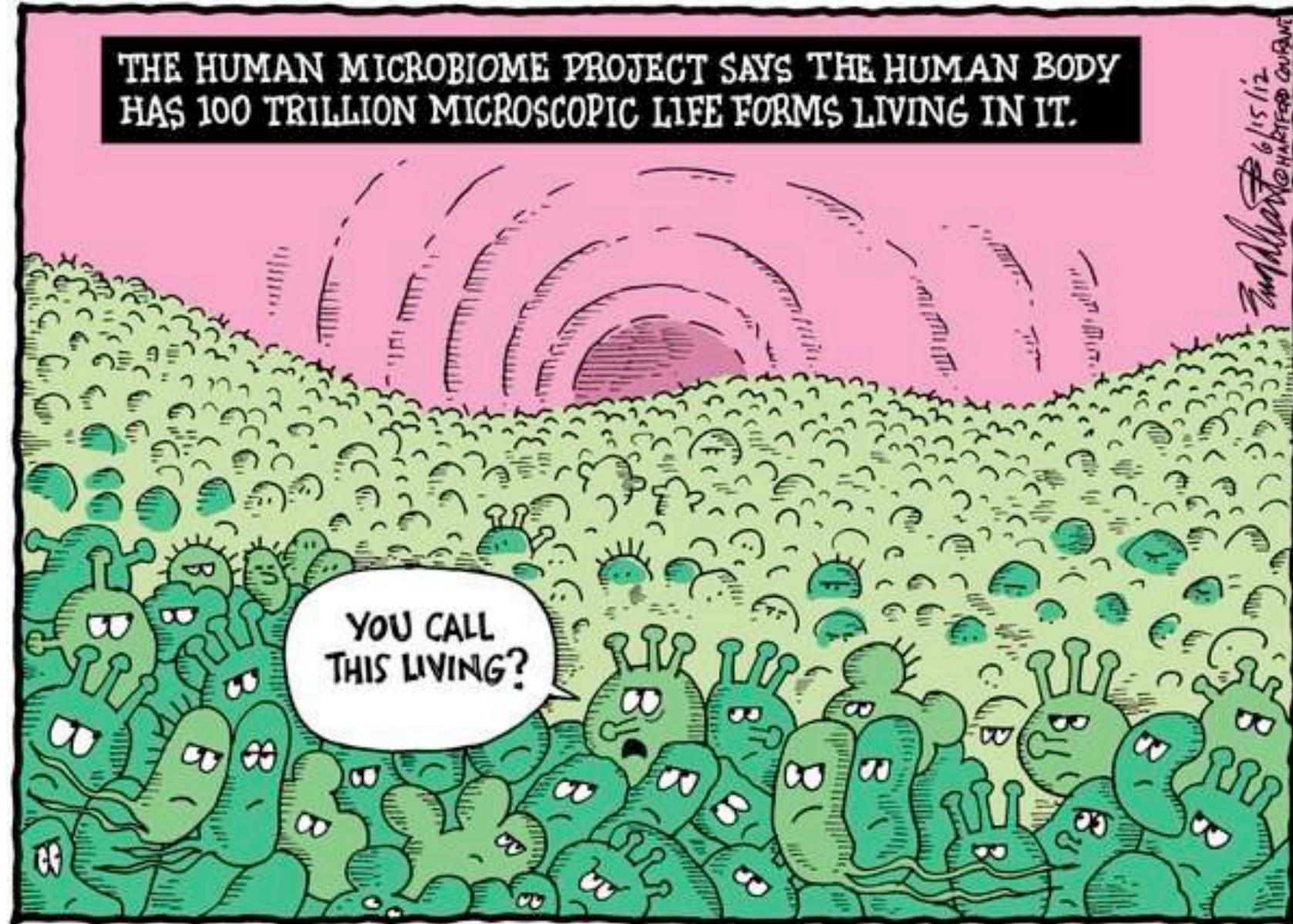
^bMendes et al. (2014) (metagenomics of rhizosphere samples).

^cTurner et al. (2013) (metatranscriptomics of rhizosphere samples).

^dBulgarelli et al. (2015) (metagenomics of rhizosphere samples).

^eQin et al. (2010) (metagenomics of gut samples).

Human gut microbiome



Human gut microbiome

Vol 464 | 4 March 2010 | doi:10.1038/nature08821

nature

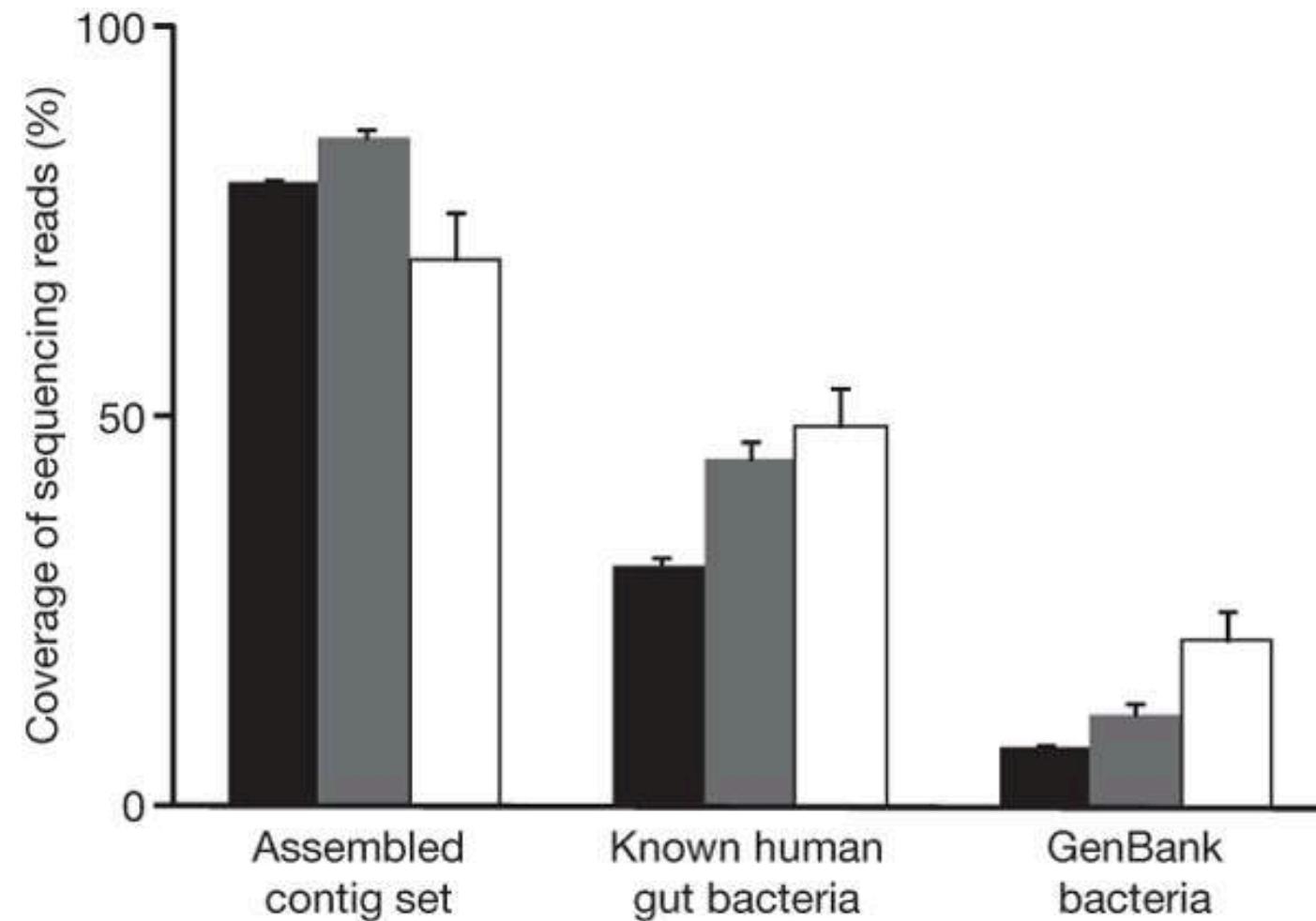
ARTICLES

A human gut microbial gene catalogue established by metagenomic sequencing

Junjie Qin^{1*}, Ruiqiang Li^{1*}, Jeroen Raes^{2,3}, Manimozhiyan Arumugam², Kristoffer Solvsten Burgdorf⁴, Chaysavanh Manichanh⁵, Trine Nielsen⁴, Nicolas Pons⁶, Florence Levenez⁶, Takuji Yamada², Daniel R. Mende², Junhua Li^{1,7}, Junming Xu¹, Shaochuan Li¹, Dongfang Li^{1,8}, Jianjun Cao¹, Bo Wang¹, Huiqing Liang¹, Huisong Zheng¹, Yinlong Xie^{1,7}, Julien Tap⁶, Patricia Lepage⁶, Marcelo Bertalan⁹, Jean-Michel Batto⁶, Torben Hansen⁴, Denis Le Paslier¹⁰, Allan Linneberg¹¹, H. Bjørn Nielsen⁹, Eric Pelletier¹⁰, Pierre Renault⁶, Thomas Sicheritz-Ponten⁹, Keith Turner¹², Hongmei Zhu¹, Chang Yu¹, Shengting Li¹, Min Jian¹, Yan Zhou¹, Yingrui Li¹, Xiuqing Zhang¹, Songgang Li¹, Nan Qin¹, Huanming Yang¹, Jian Wang¹, Søren Brunak⁹, Joel Doré⁶, Francisco Guarner⁵, Karsten Kristiansen¹³, Oluf Pedersen^{4,14}, Julian Parkhill¹², Jean Weissenbach¹⁰, MetaHIT Consortium†, Peer Bork², S. Dusko Ehrlich⁶ & Jun Wang^{1,13}

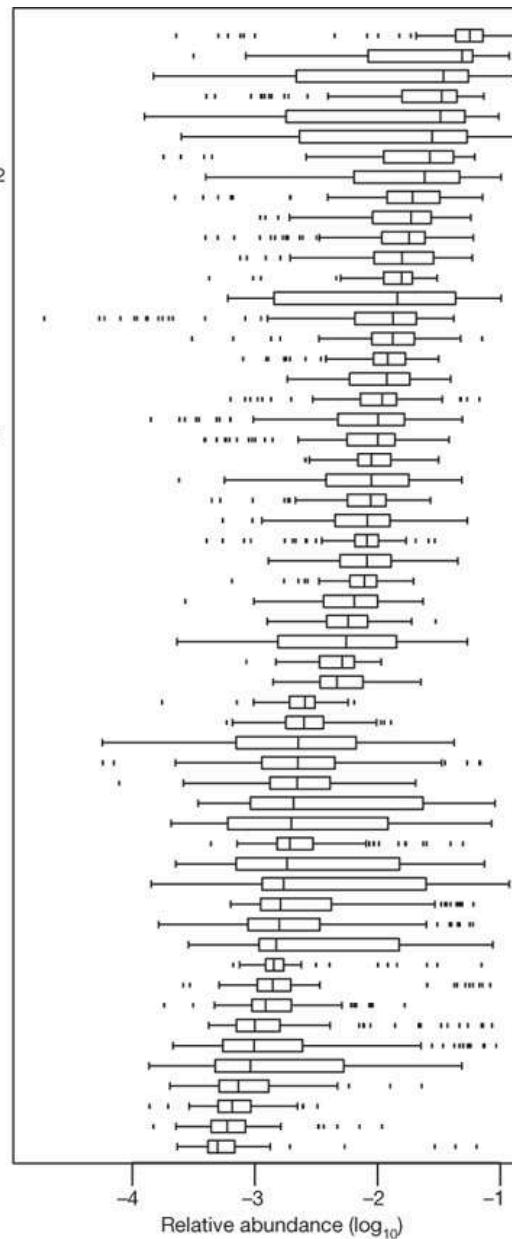
doi:10.1038/nature08821

Human gut microbiome

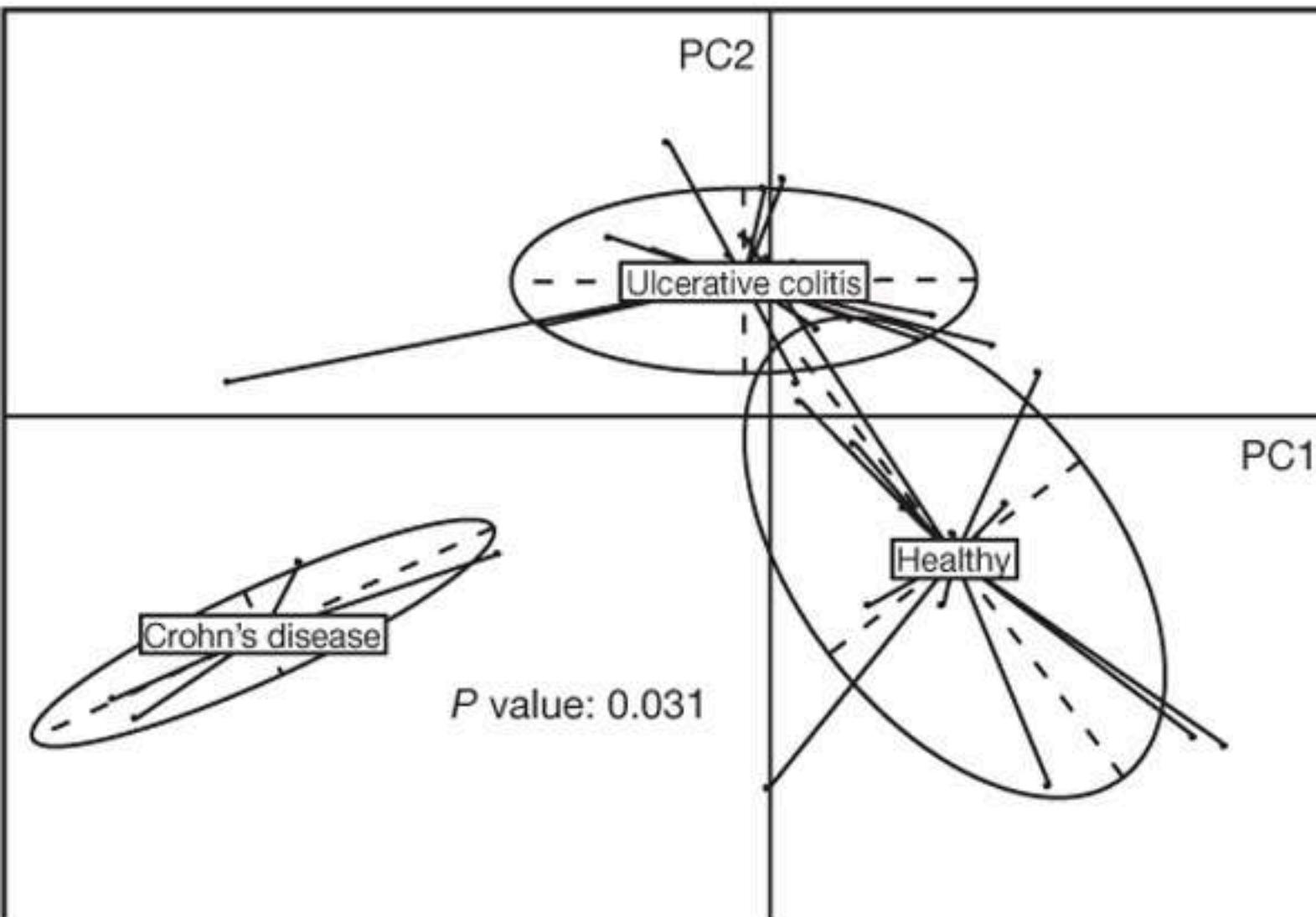


Human gut microbiome

Bacteroides uniformis
Alistipes putredinis
Parabacteroides merdae
Dorea longicatena
Ruminococcus bromii L2-63
Bacteroides caccae
Clostridium sp. SS2-1
Bacteroides thetaiotaomicron VPI-5482
Eubacterium hallii
Ruminococcus torques L2-14
Unknown sp. SS3 4
Ruminococcus sp. SR1 5
Faecalibacterium prausnitzii SL3 3
Ruminococcus lactaris
Collinsella aerofaciens
Dorea formicigenerans
Bacteroides vulgatus ATCC 8482
Roseburia intestinalis M50 1
Bacteroides sp. 2_1_7
Eubacterium siraeum 70 3
Parabacteroides distasonis ATCC 8503
Bacteroides sp. 9_1_42FAA
Bacteroides ovatus
Bacteroides sp. 4_3_47FAA
Bacteroides sp. 2_2_4
Eubacterium rectale M104 1
Bacteroides xylinisolvans XB1A
Coprococcus comes SL7 1
Bacteroides sp. D1
Bacteroides sp. D4
Eubacterium ventriosum
Bacteroides dorei
Ruminococcus obeum A2-162
Subdoligranulum variabile
Bacteroides capillosus
Streptococcus thermophilus LMD-9
Clostridium leptum
Holdemani filiformis
Bacteroides stercoris
Coprococcus eutactus
Clostridium sp. M62 1
Bacteroides eggertii
Butyrivibrio crossotus
Bacteroides finegoldii
Parabacteroides johnsonii
Clostridium sp. L2-50
Clostridium nexile
Bacteroides pectinophilus
Anaerotruncus colihominis
Ruminococcus gnavus
Bacteroides intestinalis
Bacteroides fragilis 3_1_12
Clostridium asparagiforme
Enterococcus faecalis TX0104
Clostridium scindens
Blautia hansenii



Human gut microbiome



We can check which OTUs constitute the clustering (and separation) patterns

- > Biology
- > Biomarkers

Human gut microbiome

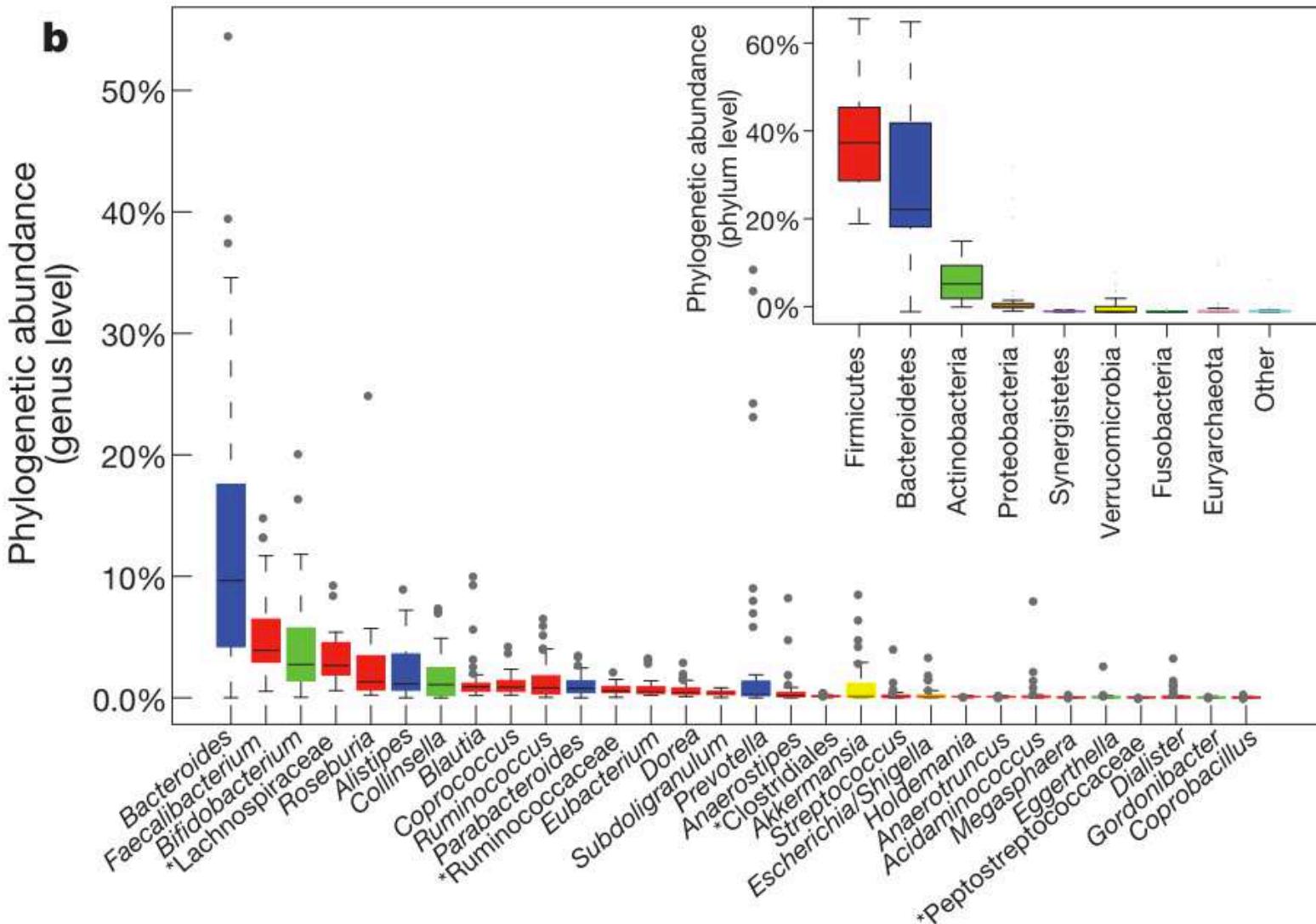




Table 1 | HMP donor samples examined by 16S and WGS

Body region	Body site	Total samples	Total 16S samples	V13 samples	V13 read depth (M)*	V35 samples	V35 read depth (M)*	Samples V13 and V35	Total WGS samples	Total read depth (G)†	Filtered reads (%)‡	Human reads (%)§	Remaining read depth (G)†	Samples 16S and WGS
Gut	Stool	352	337	193	1.4	328	2.4	184	136	1,720.7	15	1	1,450.6	124
Oral cavity	Buccal mucosa	346	330	184	1.3	314	1.7	168	107	1,438.0	9	82	136.7	91
	Hard palate	325	325	179	1.2	310	1.7	164	1	10.9	20	25	5.9	1
	Keratinized gingiva	335	329	183	1.3	319	1.7	173	6	72.3	5	47	34.4	0
	Palatine tonsils	337	332	189	1.2	315	1.9	172	6	74.8	2	80	13.5	1
	Saliva	315	310	166	0.9	292	1.5	148	5	55.7	1	91	4.2	0
	Subgingival plaque	334	328	186	1.2	314	1.8	172	7	92.1	5	79	15.3	1
	Supragingival plaque	345	331	192	1.3	316	1.9	177	115	1,500.7	15	40	674.8	101
	Throat	331	325	176	1.0	312	1.7	163	7	78.8	4	79	13.6	1
Airway	Tongue dorsum	348	332	193	1.3	320	2.0	181	122	1,620.1	15	19	1,084.3	106
	Anterior nares	316	302	169	1.0	283	1.2	150	84	1,129.9	3	96	14.3	70
Skin	Left antecubital fossa	269	269	158	0.7	221	0.5	110	0	NA	NA	NA	0	NA
	Left retroauricular crease	313	312	188	1.6	295	1.5	171	9	126.3	9	73	22.1	8
	Right antecubital fossa	274	274	158	0.7	229	0.5	113	0	NA	NA	NA	0	NA
Vagina	Right retroauricular crease	319	316	190	1.4	304	1.6	178	15	181.9	18	59	42.4	12
	Mid-vagina	145	143	91	0.6	140	1.0	88	2	22.6	0	99	0.2	0
	Posterior fornix	152	142	89	0.6	136	1.0	83	53	702.1	6	90	25.2	43
	Vaginal introitus	142	140	87	0.6	131	0.9	78	3	36.5	1	98	0.6	1
	Total	5,298	5,177	2,971	19	4,879	26.3	2,673	681	8,863.3	11	49	3,538.1	560

NCBI

6. Data submitted to NCBI Sequence Read Archives (SRA)

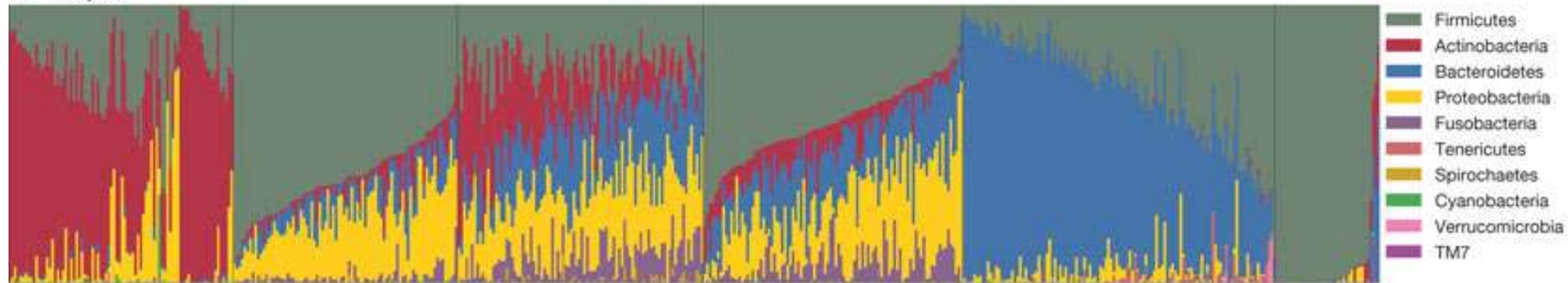
SRX: sequencing experiment

SRR: sequence run

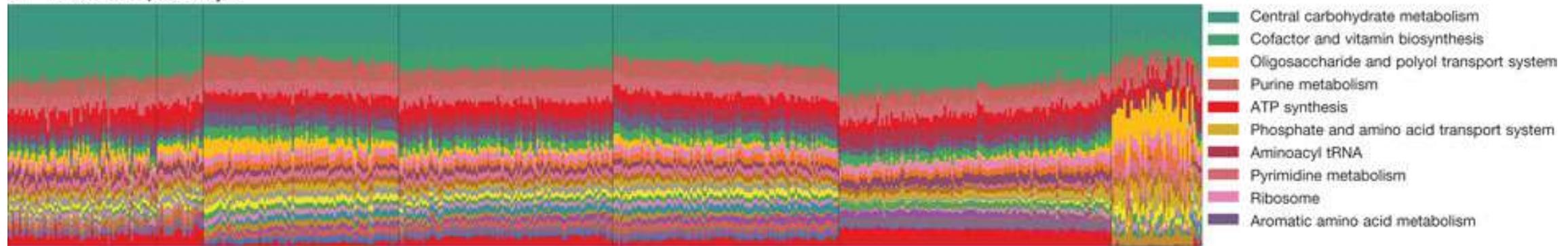
SRS: sequencing sample (maps to SN)

Human microbiome

a Phyla



b Metabolic pathways



Anterior nares

RC

Buccal mucosa

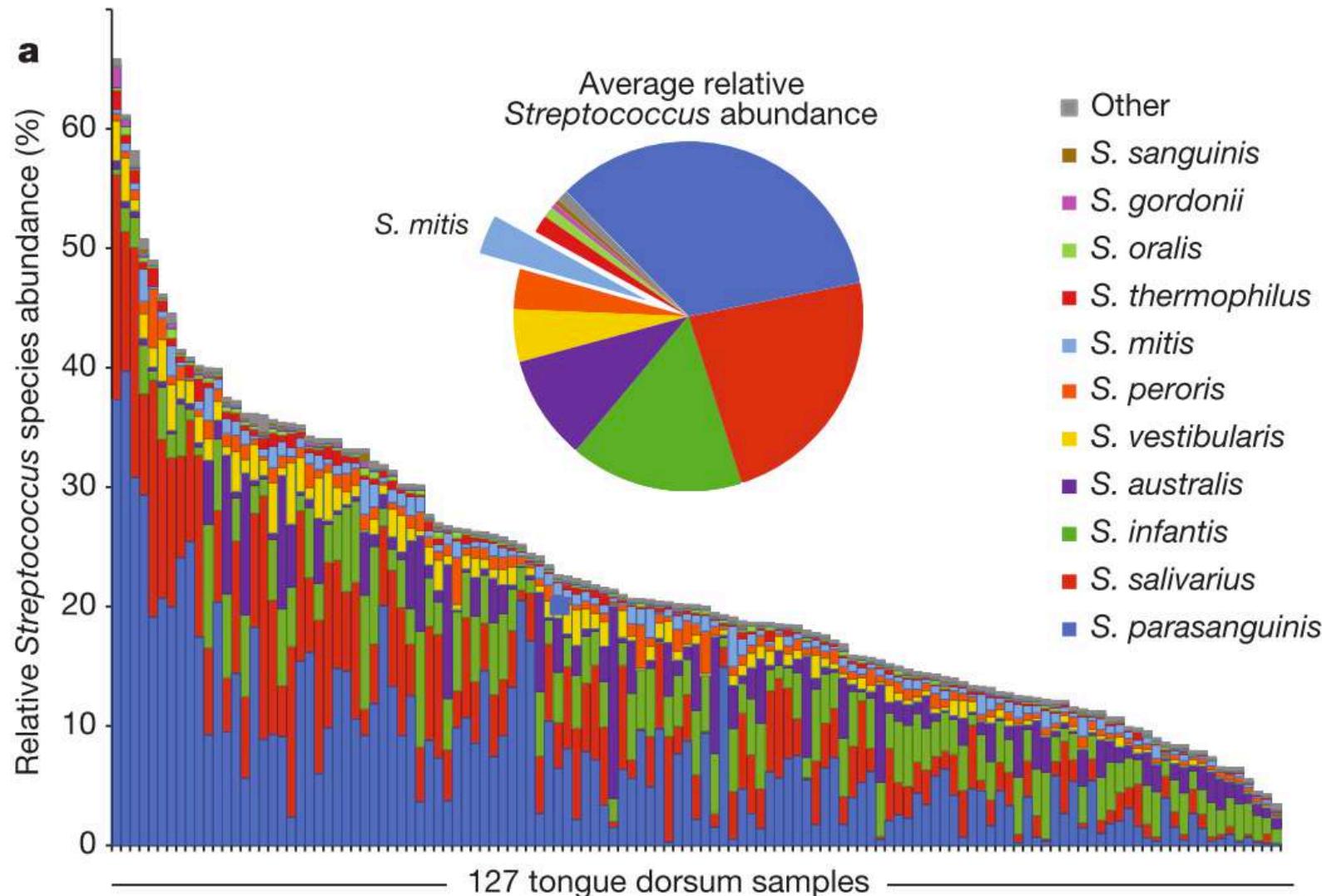
Supragingival plaque

Tongue dorsum

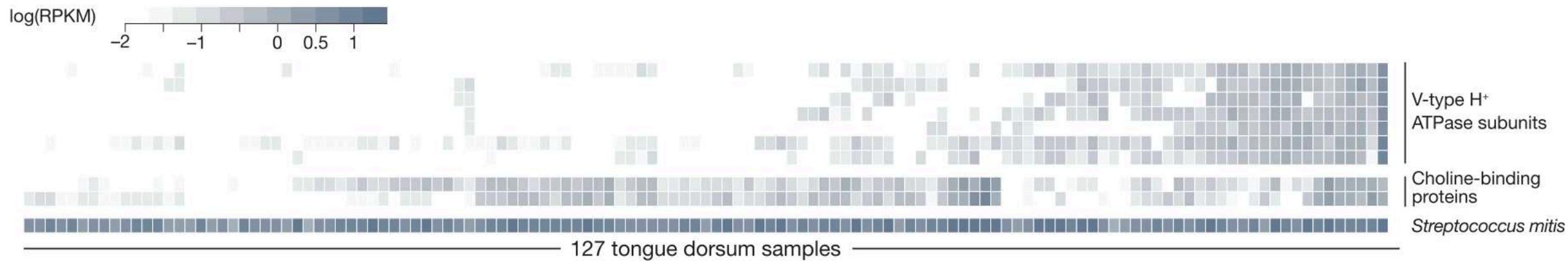
Stool

Posterior fornix

Inter-individual variation in the microbiome proved to be specific, functionally relevant and personalized



Gene loss & Structural variants are common



Skins



Genus level

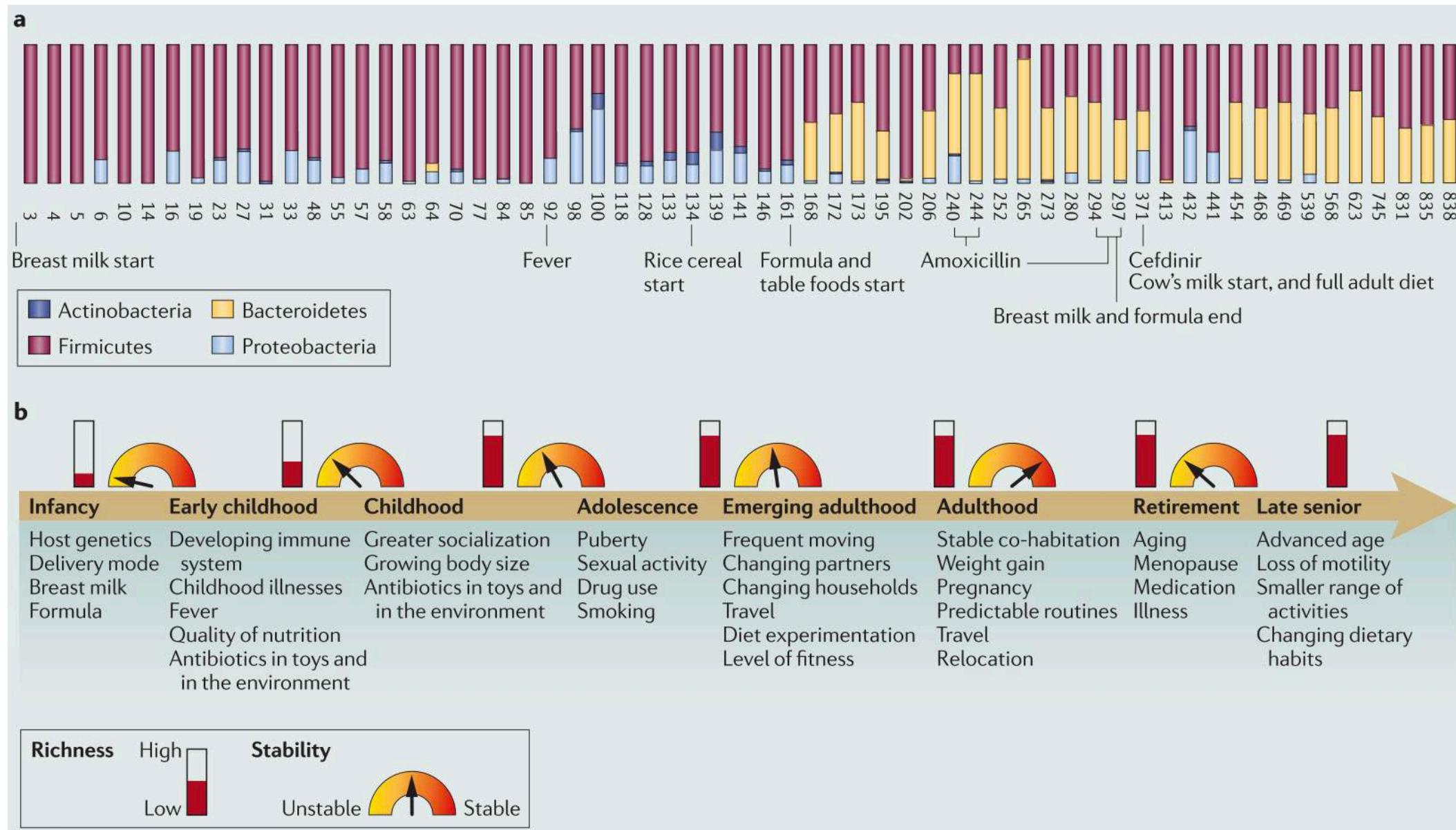
Ascomyctes

- █ Arthrodermataceae
- █ Aspergillus
- █ Candida
- █ Chrysosporium
- █ Epicoccum
- █ Leptosphaerulina
- █ Penicillium
- █ Phoma
- █ Saccharomyces

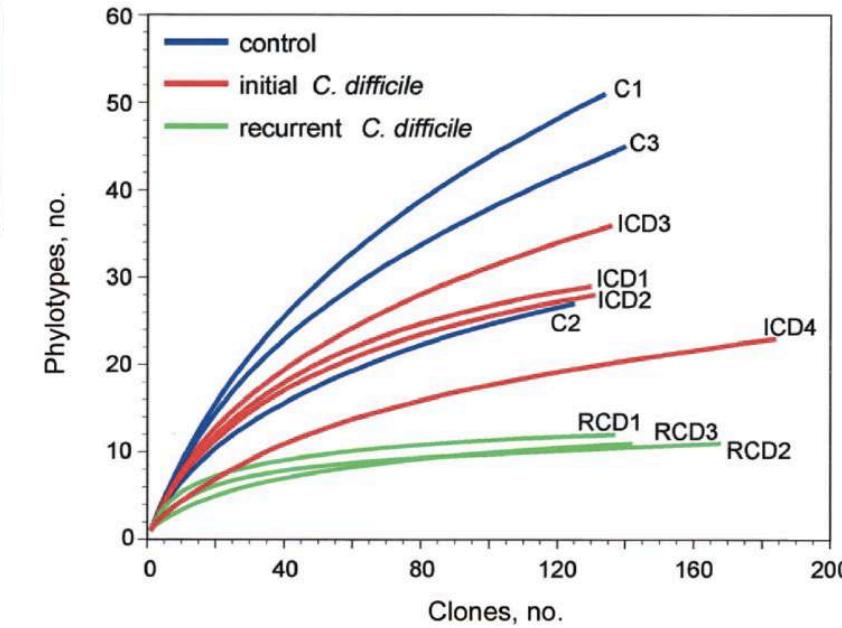
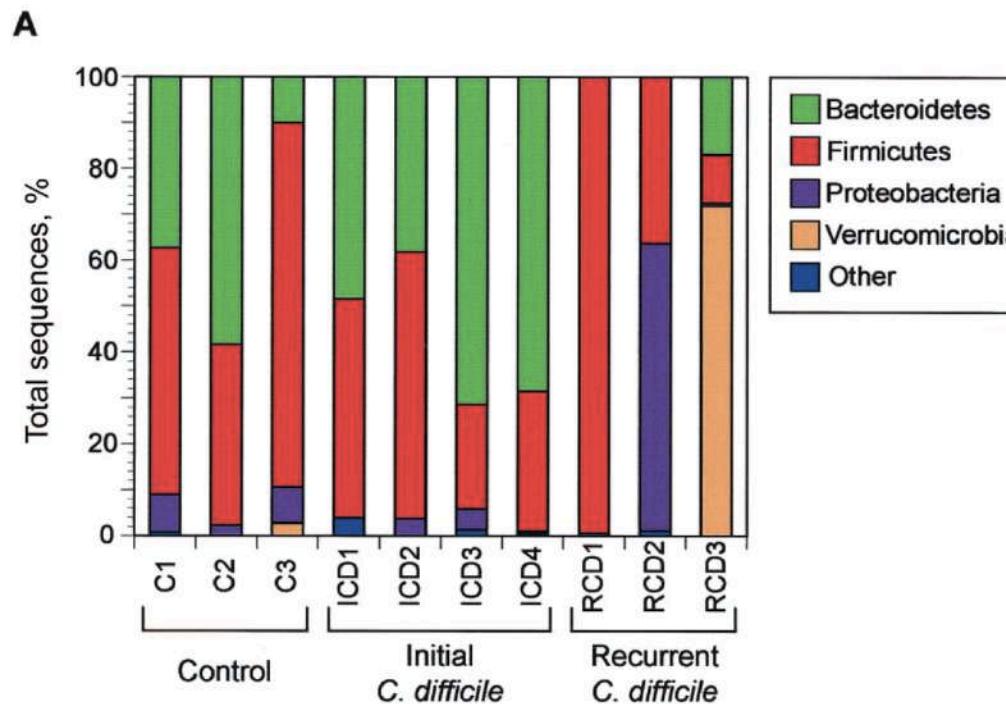
Basidiomycetes

- █ Cryptococcus
- █ Malassezia
- █ Rhodotorula
- █ Ustilago
- █ Others (<1%)

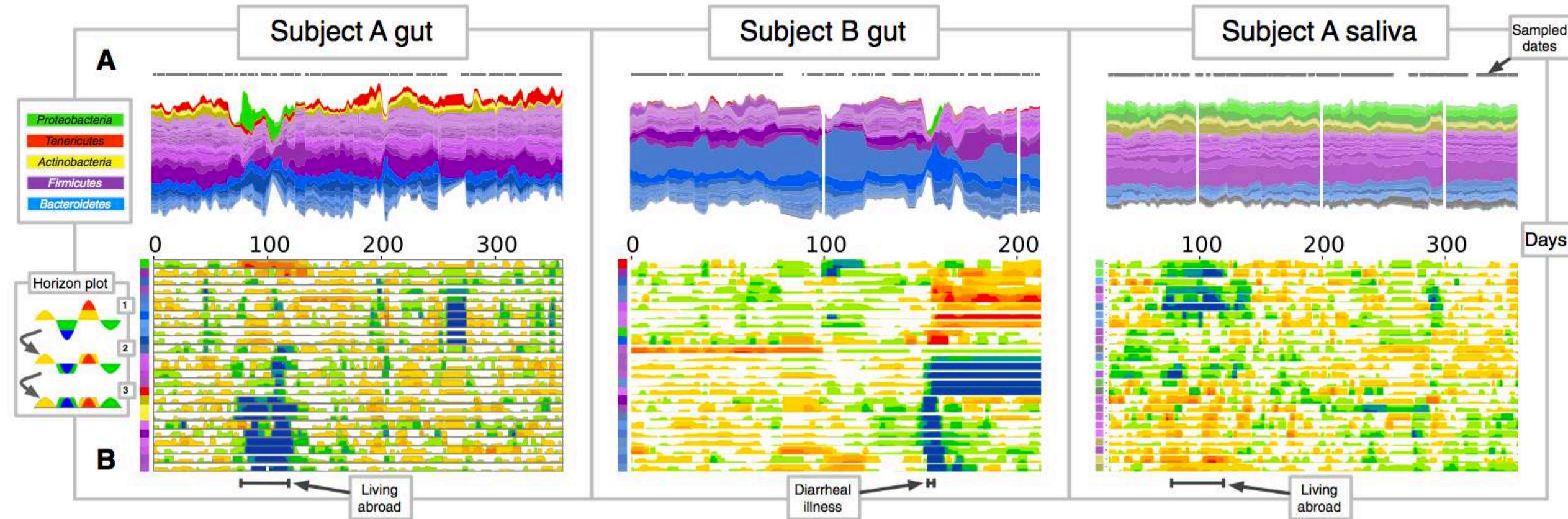
The gut microbiome during life



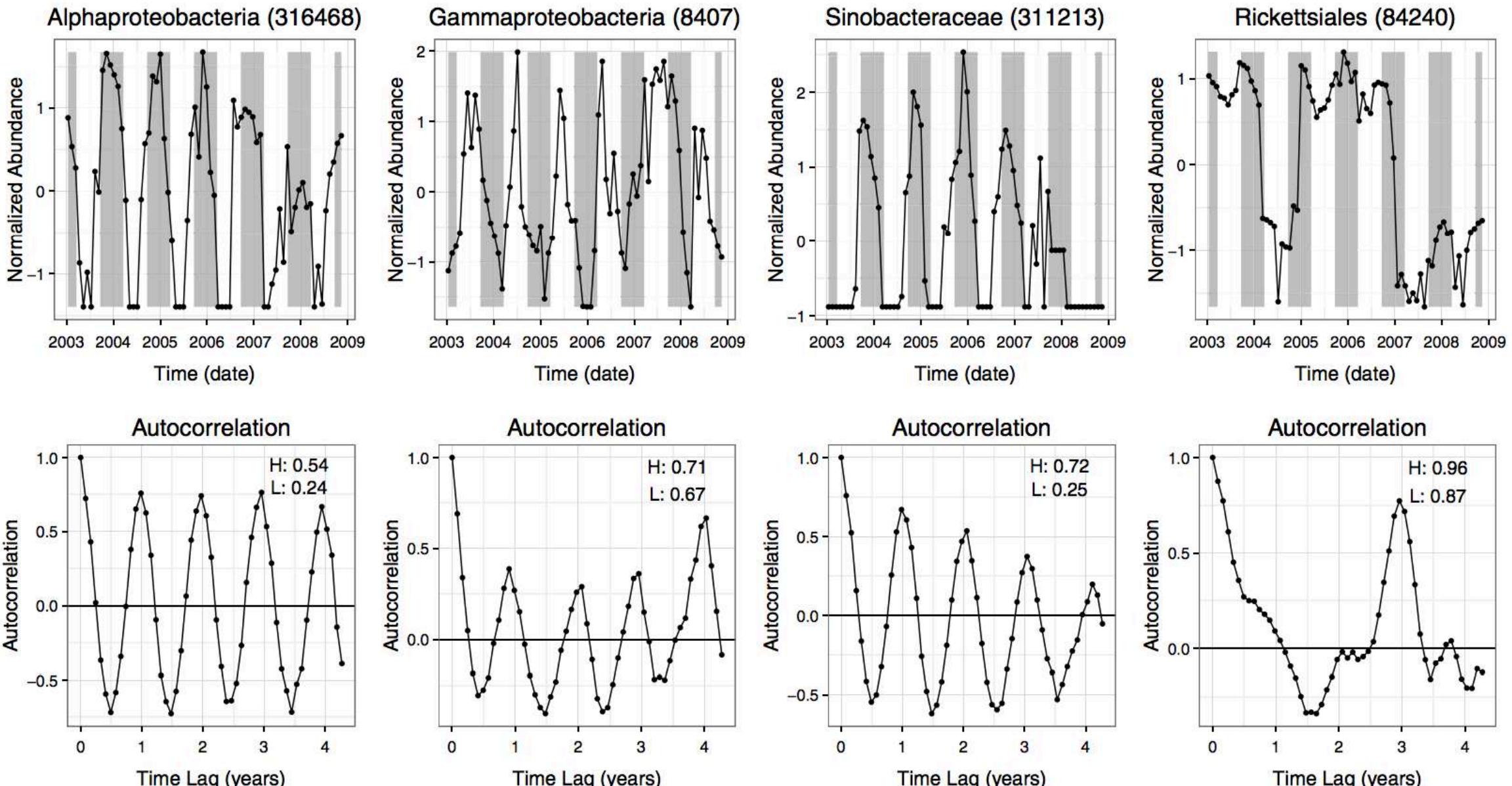
Decreased diversity with *Clostridium difficile* – associated diarrhea



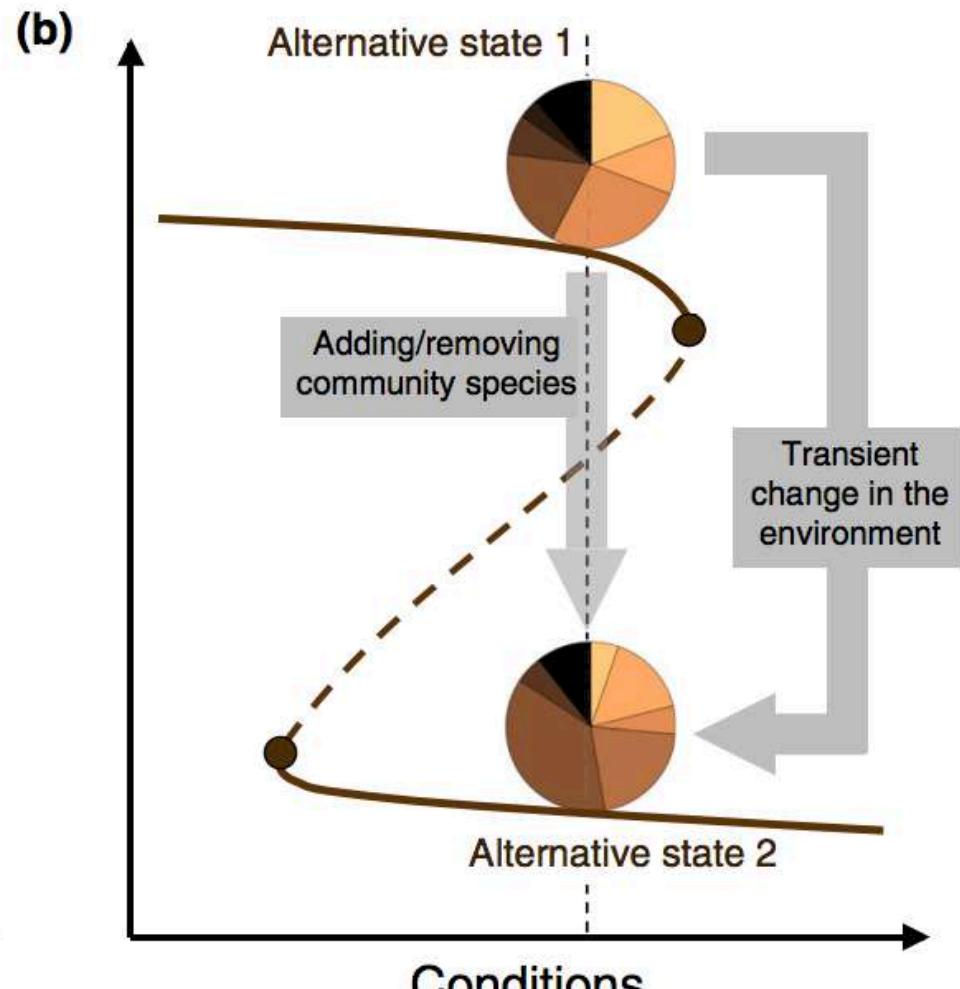
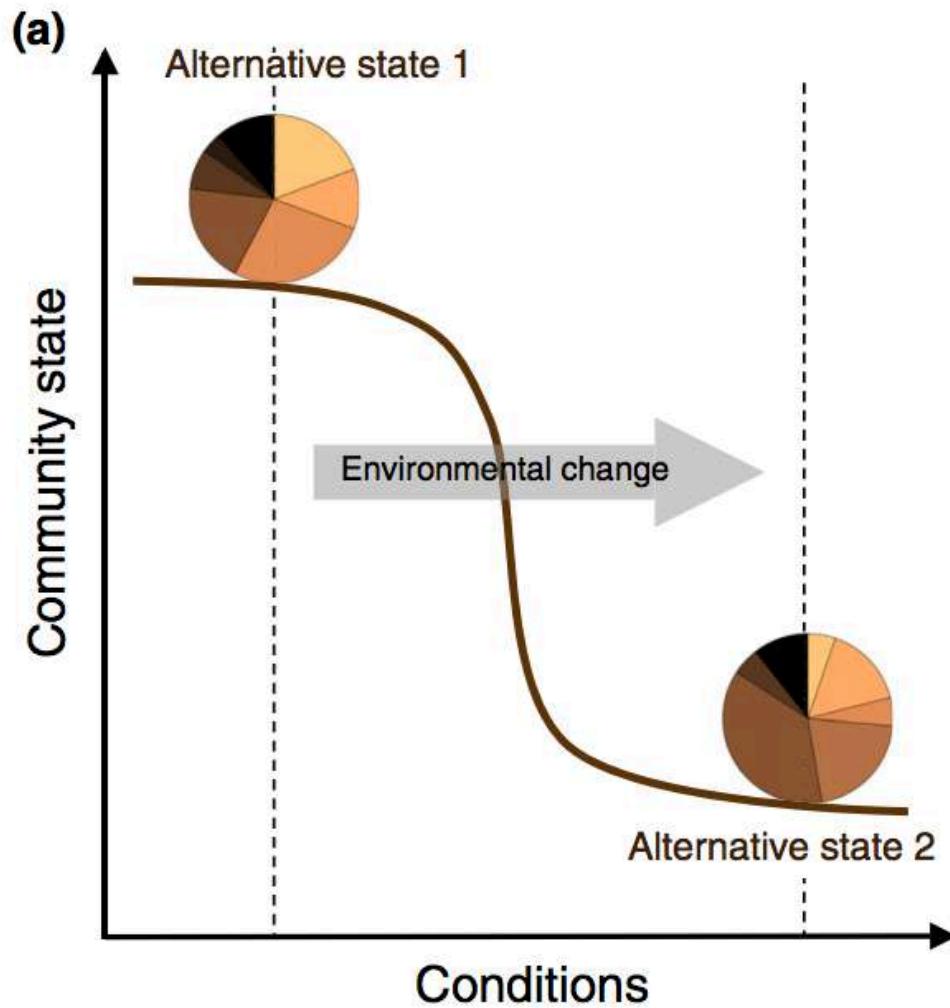
Tracking microbiome on a daily scale



Tracking microbiome spanning 6 years



Tracking microbiome on a daily scale



Current Opinion in Microbiology

Question: What community gets reset and what don't?

A. Shade, J.S. Read, N.D. Youngblut, N. Fierer, R. Knight, T.K. Kratz, N.R. Lottig, E.E. Roden, E.H. Stanley, J. Stombaugh, et al.

Lake microbial communities are resilient after a whole-ecosystem disturbance **Yes**
ISME J, 6 (2012), pp. 2153–2167

L. Dethlefsen, D.A. Relman

Incomplete recovery and individualized responses of the human distal gut microbiota to repeated antibiotic perturbation

Proc Natl Acad Sci U S A, 108 (2011), pp. 4554–4561

No

L.A. David, A.C. Materna, J. Friedman, M.I. Campos-Baptista, M.C. Blackburn, A. Perrotta, S.E. Erdman, E.J. Alm

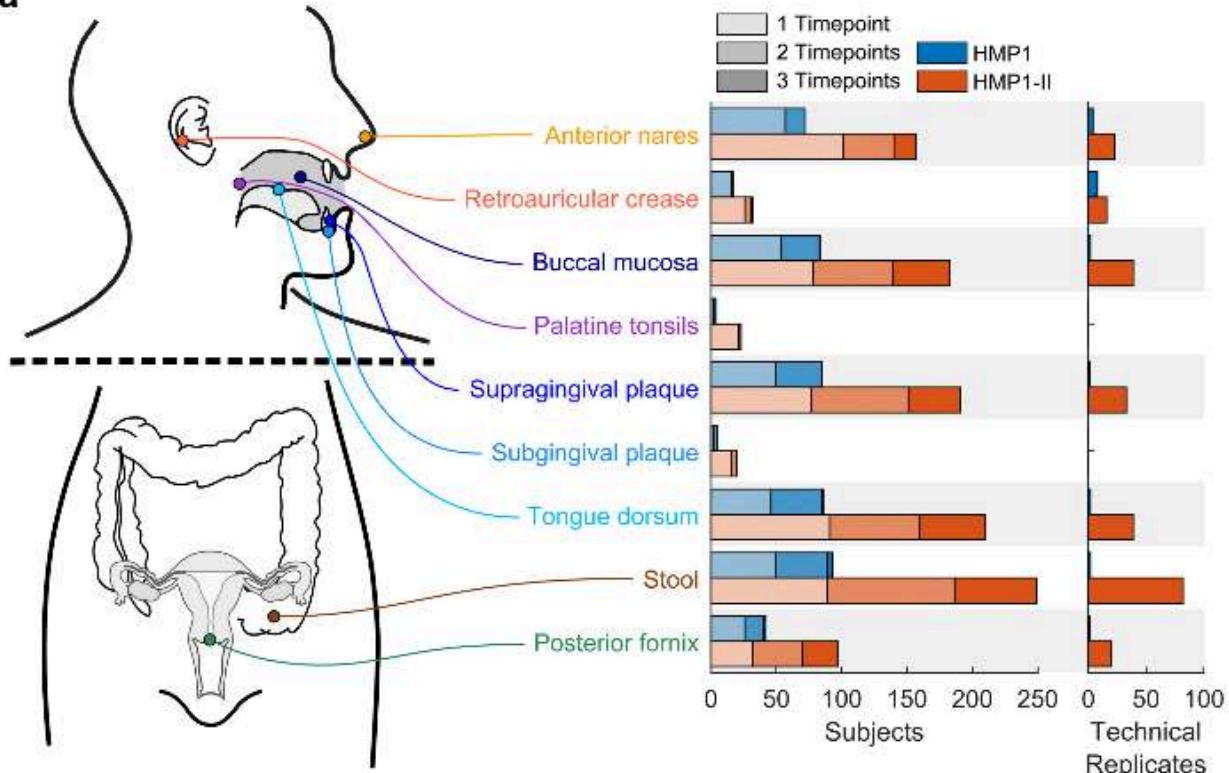
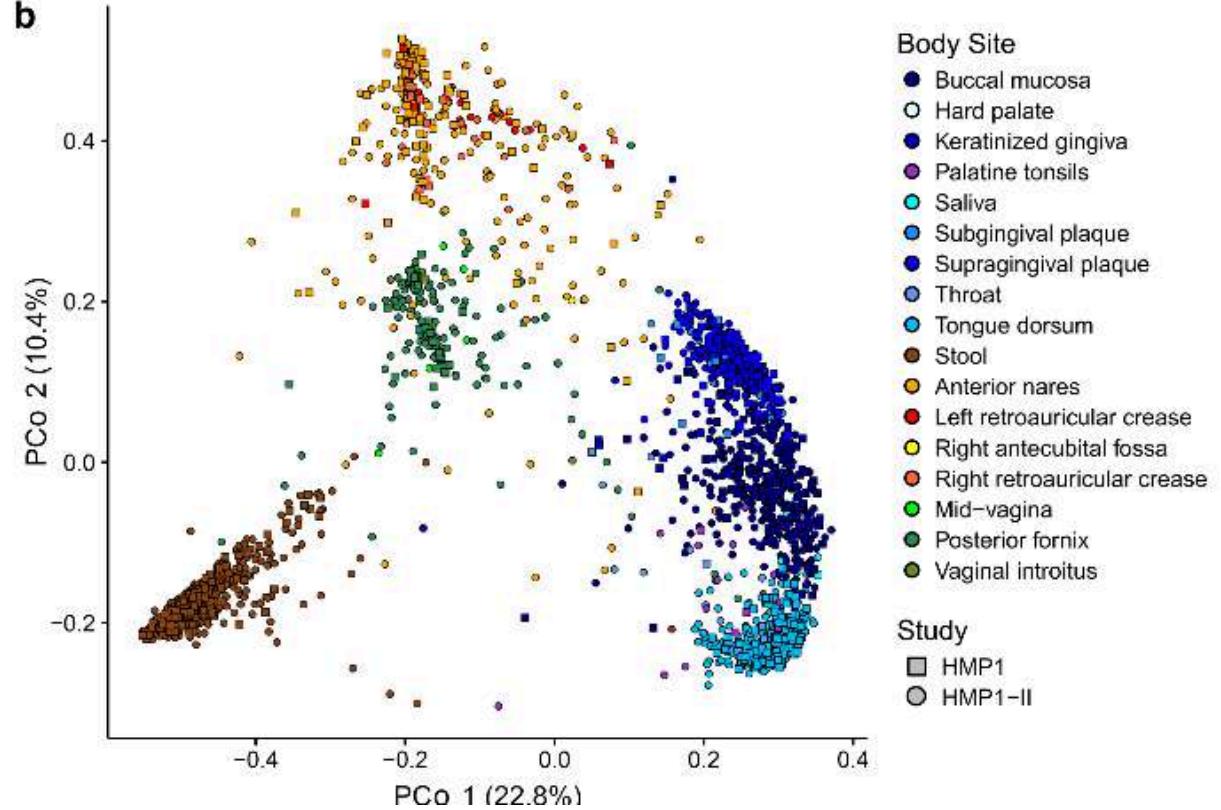
Host lifestyle affects human microbiota on daily timescales
Genome Biol, 15 (2014), p. R89

Yes and No

Strains, functions and dynamics in the expanded Human Microbiome Project

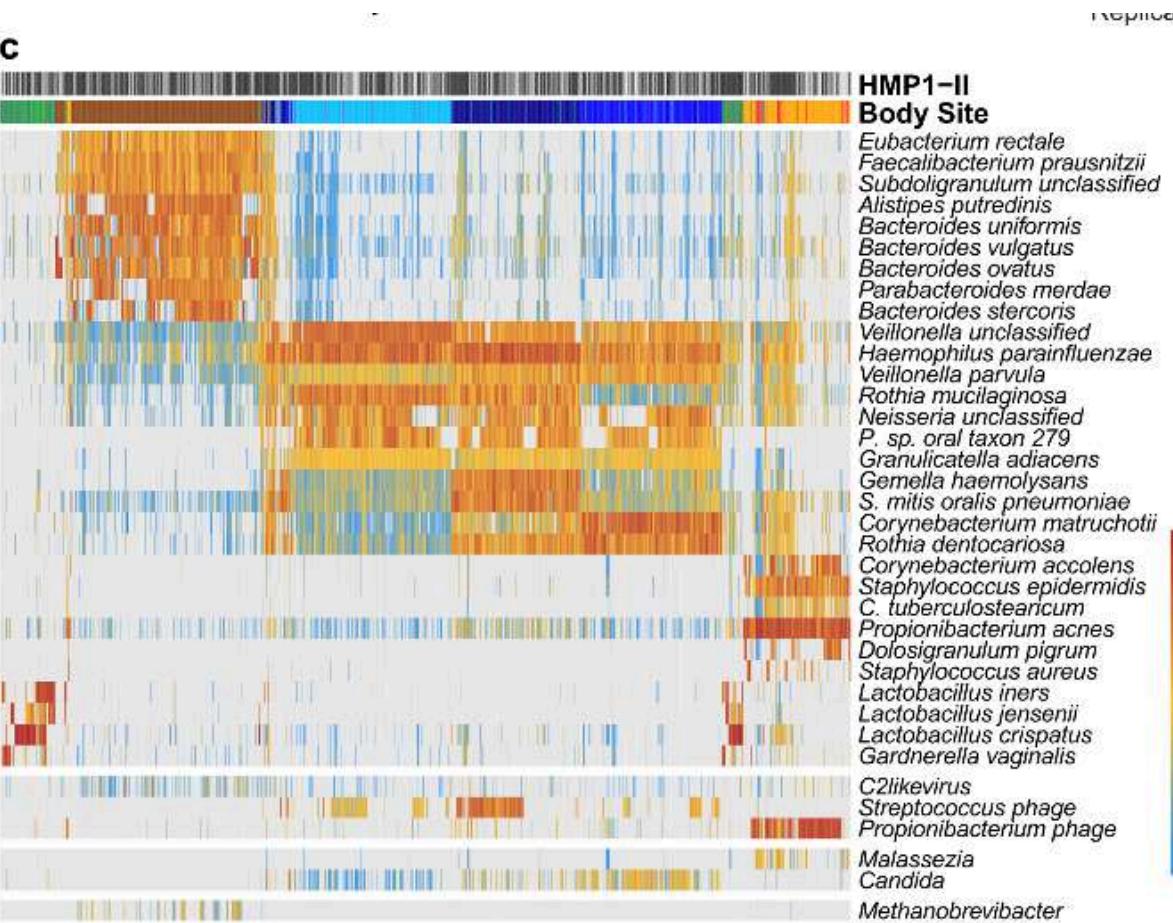
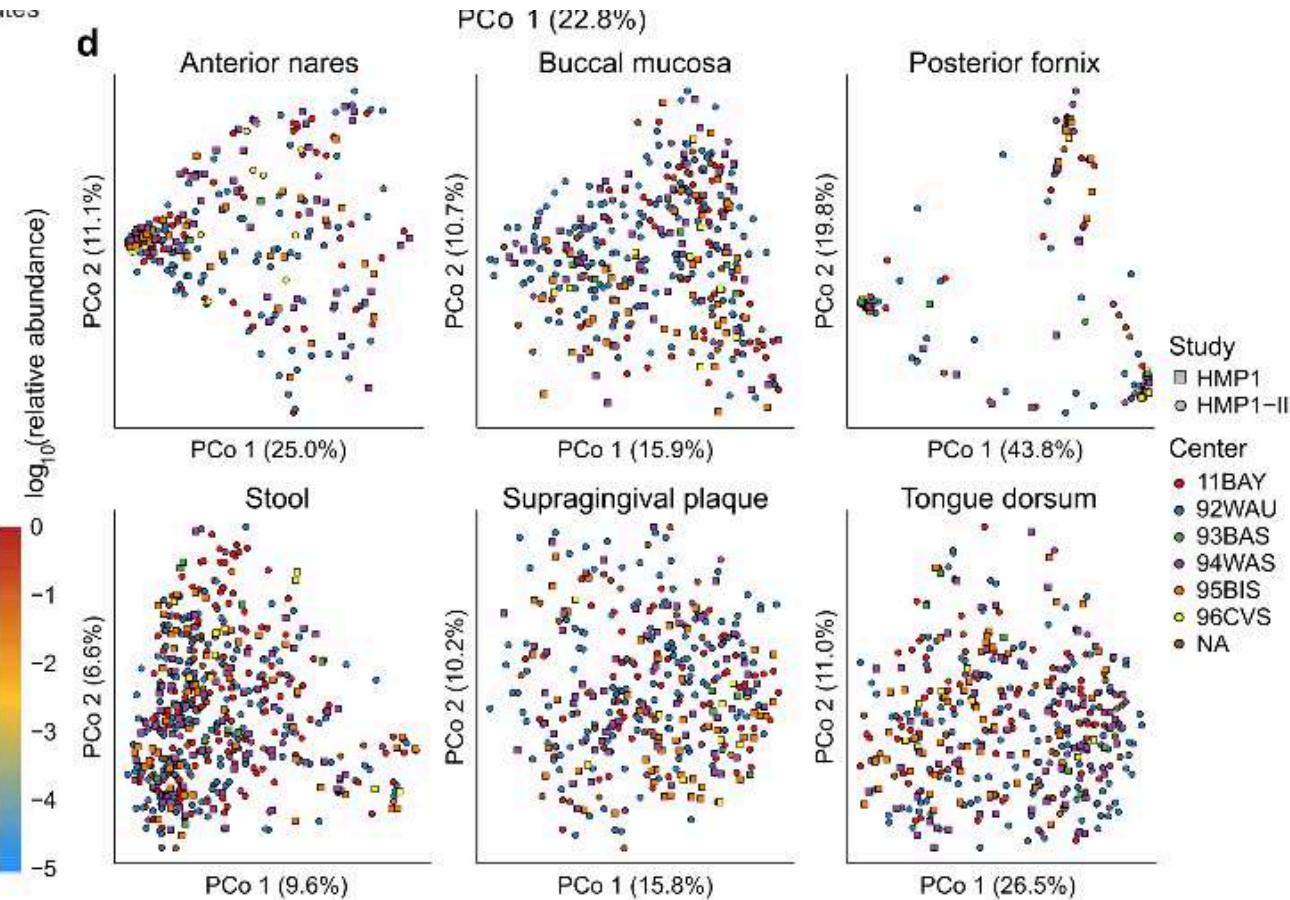
Jason Lloyd-Price^{1,2*}, Anup Mahurkar^{3*}, Gholamali Rahnavard^{1,2}, Jonathan Crabtree³, Joshua Orvis³, A. Brantley Hall², Arthur Brady³, Heather H. Creasy³, Carrie McCracken³, Michelle G. Giglio³, Daniel McDonald⁴, Eric A. Franzosa^{1,2}, Rob Knight^{4,5}, Owen White³ & Curtis Huttenhower^{1,2}

The characterization of baseline microbial and functional diversity in the human microbiome has enabled studies of microbiome-related disease, diversity, biogeography, and molecular function. The National Institutes of Health Human Microbiome Project has provided one of the broadest such characterizations so far. Here we introduce a second wave of data from the study, comprising 1,631 new metagenomes (2,355 total) targeting diverse body sites with multiple time points in 265 individuals. We applied updated profiling and assembly methods to provide new characterizations of microbiome personalization. Strain identification revealed subspecies clades specific to body sites; it also quantified species with phylogenetic diversity under-represented in isolate genomes. Body-wide functional profiling classified pathways into universal, human-enriched, and body site-enriched subsets. Finally, temporal analysis decomposed microbial variation into rapidly variable, moderately variable, and stable subsets. This study furthers our knowledge of baseline human microbial diversity and enables an understanding of personalized microbiome function and dynamics.

a**b**

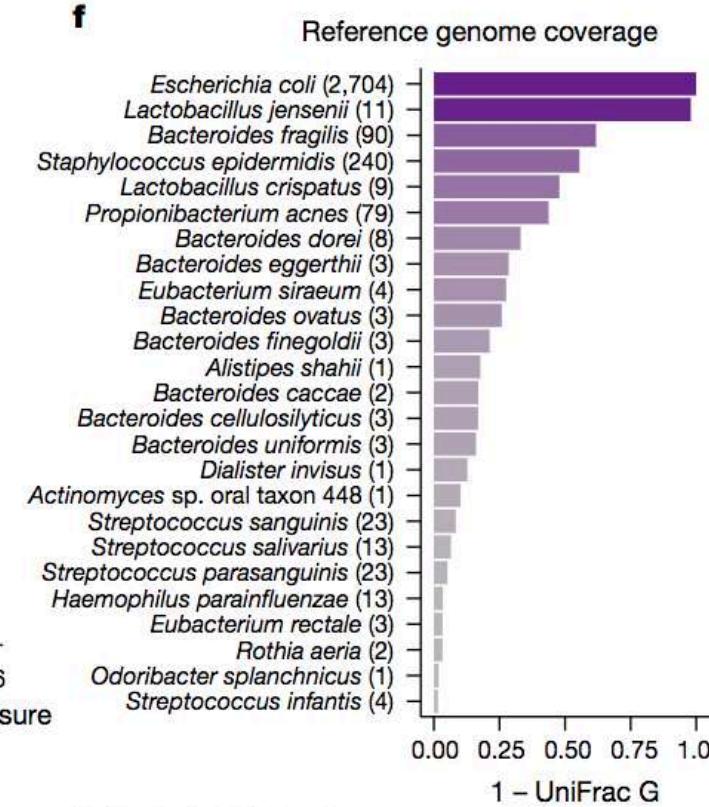
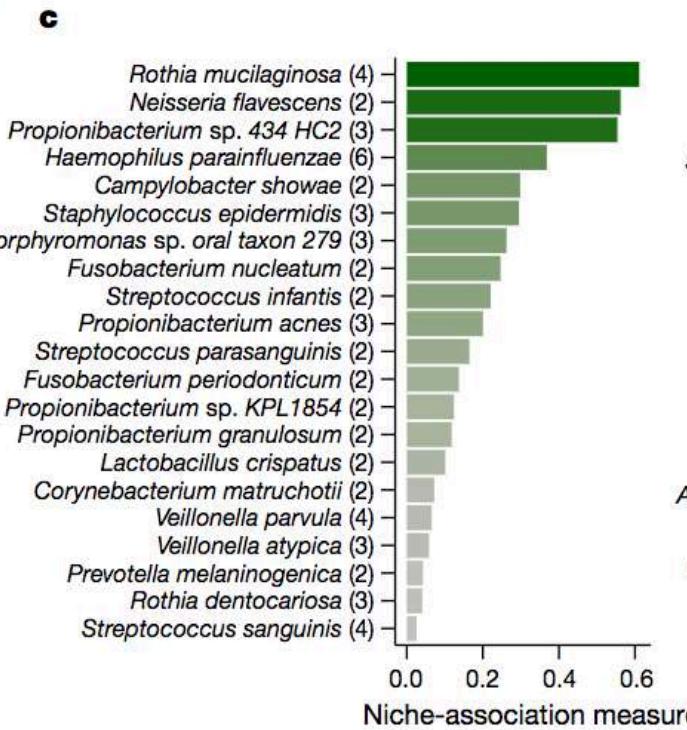
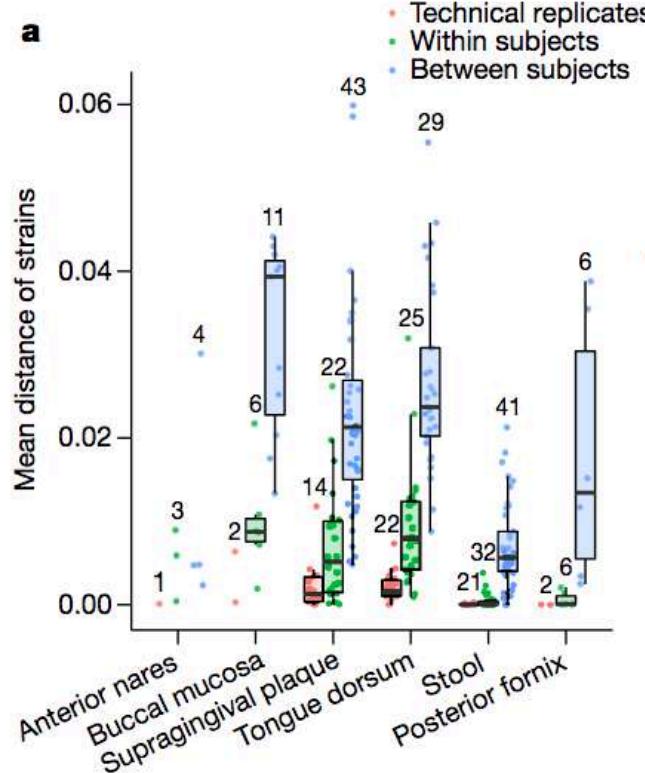
Extended Data Figure 1 | Extended body-wide metagenomic taxonomic profiles in HMP1-II. **a**, The combined HMP1-II datasets include a total of 2,355 metagenomes (724 previously published and 1,631 new, including 252 technical replicates). These span the project's six targeted body sites (anterior nares, buccal mucosa, supragingival plaque, tongue dorsum, stool, and posterior fornix) in addition to at least 20 samples each from 3 additional sites, of the 18 total sampled sites: retroauricular crease, palatine tonsils, and subgingival plaque. Metagenomes are now available for at least one body site for a total of 265 individuals. **b**, PCoA using Bray–Curtis

distances among all microbes at the species level. **c**, Relative abundances of the most prevalent and abundant microbes (bacterial, viral, eukaryotic, and archaeal) among all body sites, as profiled by MetaPhlAn2²⁰. Prevalent eukaryotic microbes are shown at the genus level. **d**, Taxonomic profiles do not vary more between sequencing centres, batches, or clinical centres than they do among individuals within body sites. Ordinations show Bray–Curtis principal coordinates of species-level abundances at each body site. Within-site ecological structure is as expected¹, with no divergence associated with technical variables along the first two ordination axes.

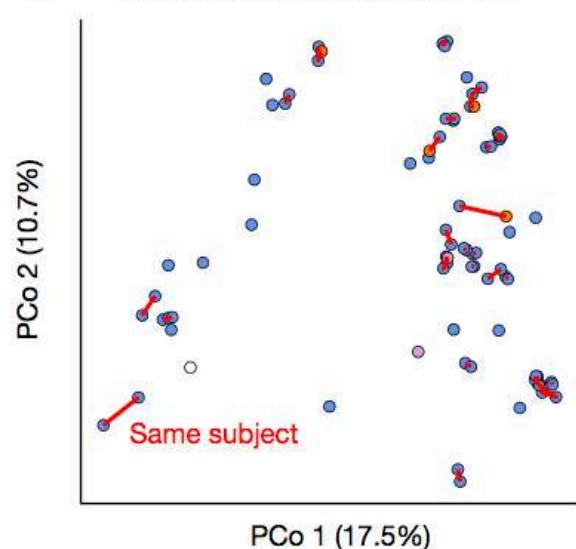
c**d**

Extended Data Figure 1 | Extended body-wide metagenomic taxonomic profiles in HMP1-II. **a**, The combined HMP1-II datasets include a total of 2,355 metagenomes (724 previously published and 1,631 new, including 252 technical replicates). These span the project's six targeted body sites (anterior nares, buccal mucosa, supragingival plaque, tongue dorsum, stool, and posterior fornix) in addition to at least 20 samples each from 3 additional sites, of the 18 total sampled sites: retroauricular crease, palatine tonsils, and subgingival plaque. Metagenomes are now available for at least one body site for a total of 265 individuals. **b**, PCoA using Bray–Curtis

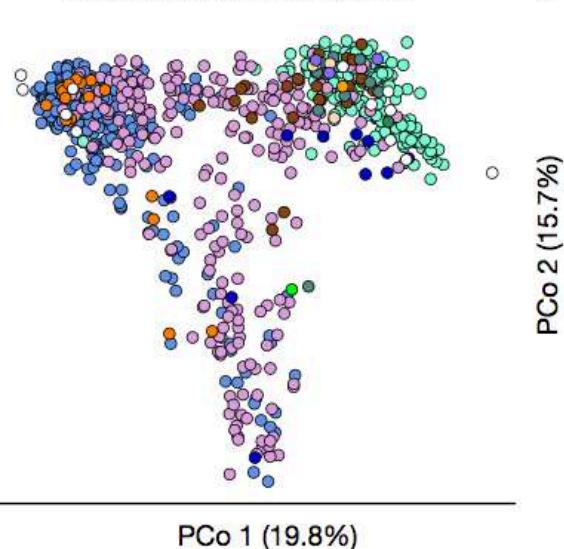
distances among all microbes at the species level. **c**, Relative abundances of the most prevalent and abundant microbes (bacterial, viral, eukaryotic, and archaeal) among all body sites, as profiled by MetaPhlAn2²⁰. Prevalent eukaryotic microbes are shown at the genus level. **d**, Taxonomic profiles do not vary more between sequencing centres, batches, or clinical centres than they do among individuals within body sites. Ordinations show Bray–Curtis principal coordinates of species-level abundances at each body site. Within-site ecological structure is as expected¹, with no divergence associated with technical variables along the first two ordination axes.



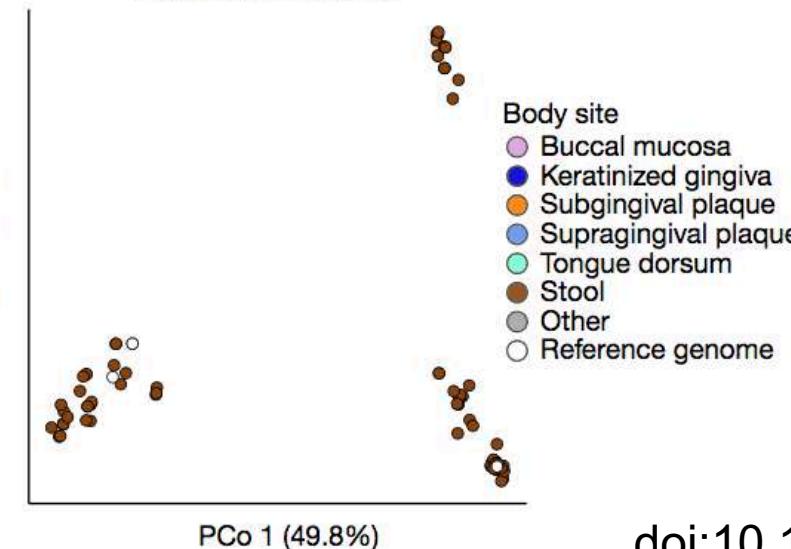
b *Actinomyces* sp. oral taxon 448

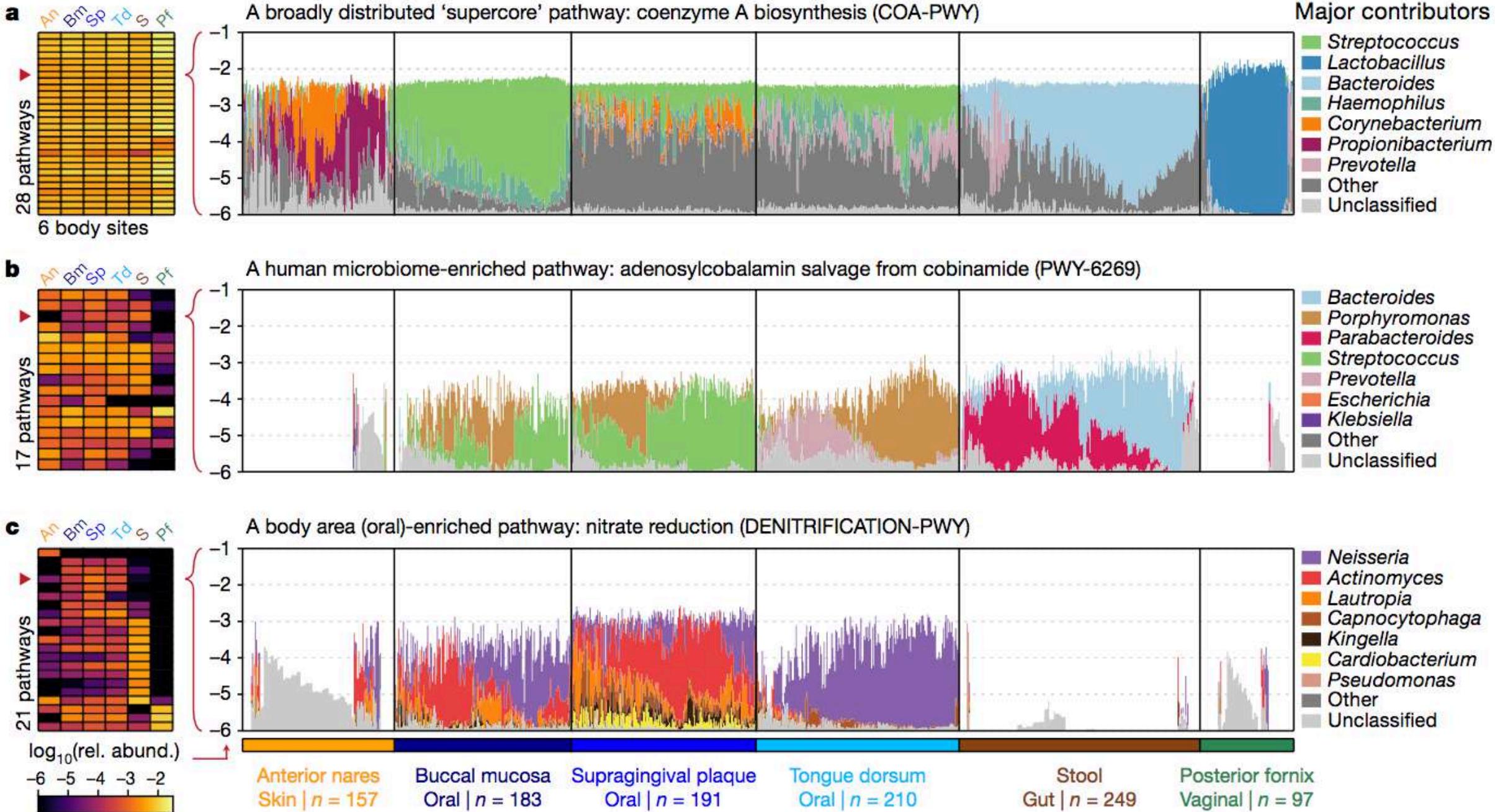


d *Haemophilus parainfluenzae*



e *Eubacterium siraeum*





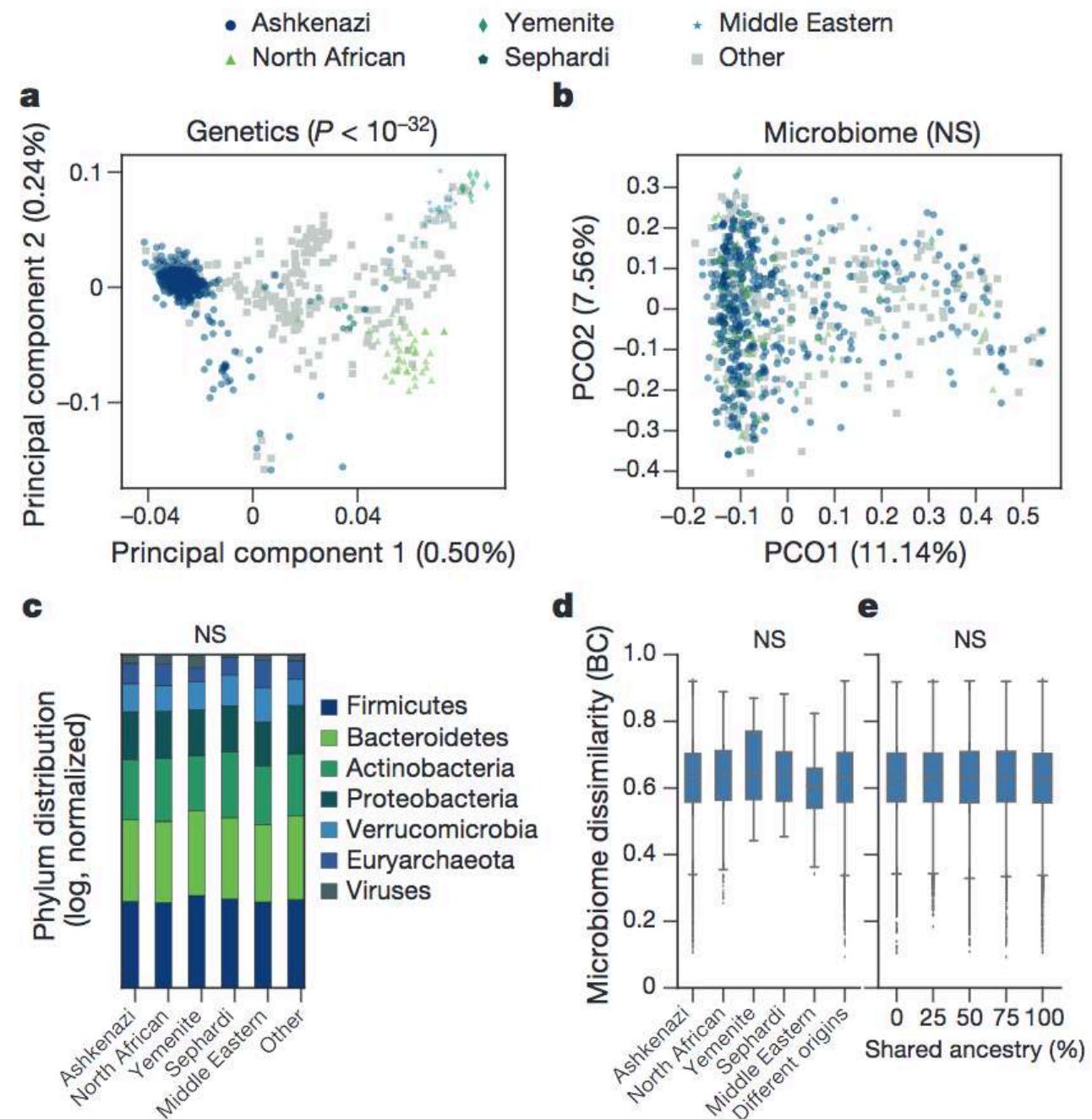
Environment dominates over host genetics in shaping human gut microbiota

Daphna Rothschild^{1,2*}, Omer Weissbrod^{1,2*}, Elad Barkan^{1,2*}, Alexander Kurilshikov³, Tal Korem^{1,2}, David Zeevi^{1,2}, Paul I. Costea^{1,2}, Anastasia Godneva^{1,2}, Iris N. Kalka^{1,2}, Noam Bar^{1,2}, Smadar Shilo^{1,2}, Dar Lador^{1,2}, Arnau Vich Vila^{3,4}, Niv Zmora^{5,6,7}, Meirav Pevsner-Fischer⁵, David Israeli⁸, Noa Kosower^{1,2}, Gal Malka^{1,2}, Bat Chen Wolf^{1,2}, Tali Avnit-Sagi^{1,2}, Maya Lotan-Pompan^{1,2}, Adina Weinberger^{1,2}, Zamir Halpern^{7,9}, Shai Carmi¹⁰, Jingyuan Fu^{3,11}, Cisca Wijmenga^{3,12}, Alexandra Zhernakova³, Eran Elinav⁵§ & Eran Segal^{1,2}§

Human gut microbiome composition is shaped by multiple factors but the relative contribution of host genetics remains elusive. Here we examine genotype and microbiome data from 1,046 healthy individuals with several distinct ancestral origins who share a relatively common environment, and demonstrate that the gut microbiome is not significantly associated with genetic ancestry, and that host genetics have a minor role in determining microbiome composition. We show that, by contrast, there are significant similarities in the compositions of the microbiomes of genetically unrelated individuals who share a household, and that over 20% of the inter-person microbiome variability is associated with factors related to diet, drugs and anthropometric measurements. We further demonstrate that microbiome data significantly improve the prediction accuracy for many human traits, such as glucose and obesity measures, compared to models that use only host genetic and environmental data. These results suggest that microbiome alterations aimed at improving clinical outcomes may be carried out across diverse genetic backgrounds.

- 1,046 healthy Israeli adults
- 16S rRNA + metagenomics
- Genotyping 712,540 SNPs
- Questionnaires

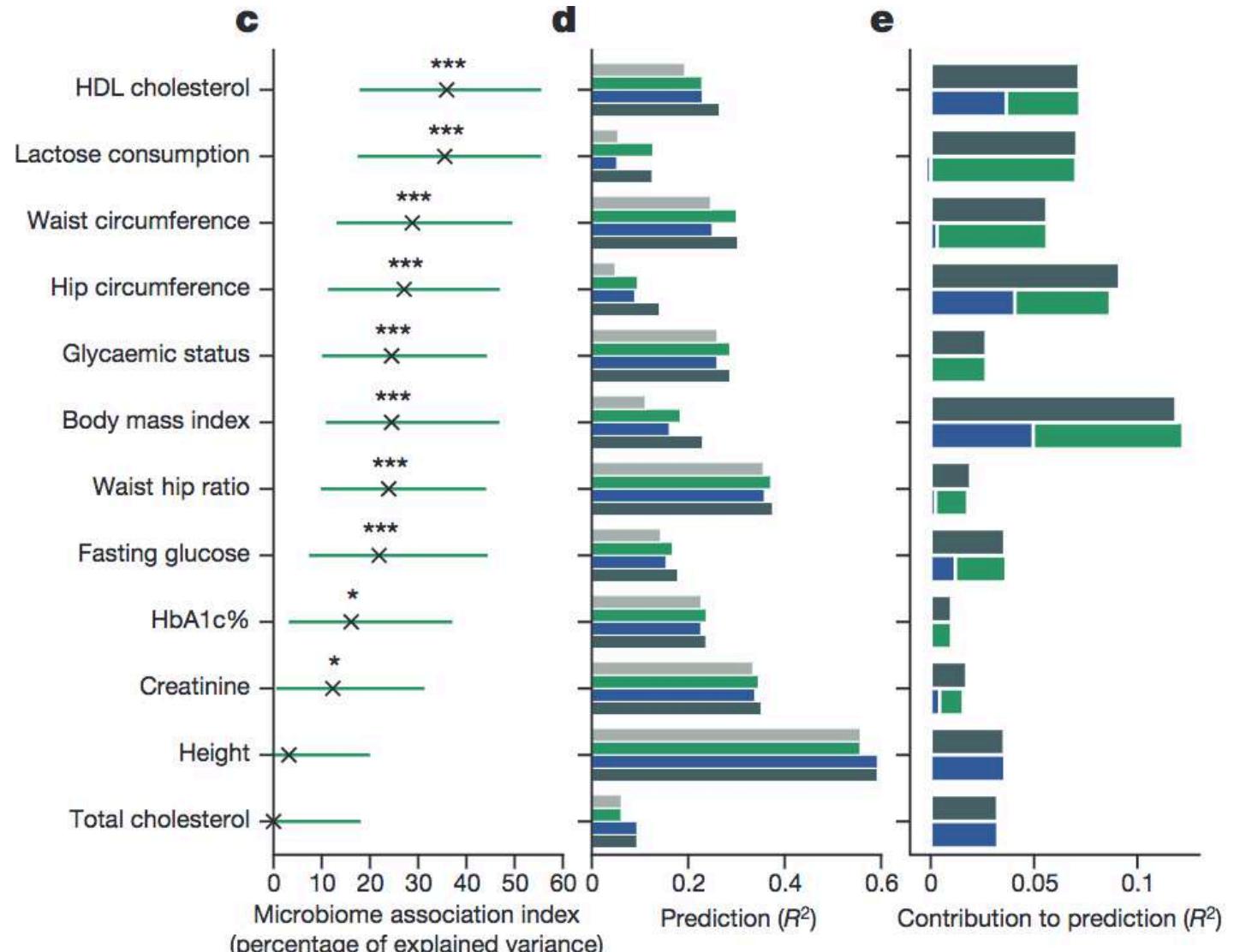
Figure 1 | Genetic ancestry is not significantly associated with microbiome composition. **a**, Genetic principal components are strongly associated with self-reported ancestry, with Ashkenazi ($n=345$), North African ($n=42$), Middle Eastern ($n=24$), Sephardi ($n=10$), Yemenite ($n=8$) and admixed/other (other) ($n=286$) ancestries ($P < 10^{-32}$; Kruskal–Wallis). **b**, As in **a**, but for microbiome principal coordinate analysis ($P > 0.08$; Kruskal–Wallis). **c**, The distribution of average phylum abundance among 582 non-admixed individuals (in log scale, normalized to sum to 1.0) is not associated with ancestry ($P > 0.05$; Kruskal–Wallis). NS, not significant. **d**, Box plots of Bray–Curtis (BC) dissimilarities across all pairs of 737 individuals for whom the ancestries of all grandparents are known, demonstrating that microbiome composition is not associated with ancestry ($P > 0.06$; Kruskal–Wallis test for the top five Bray–Curtis PCOs). $n = 105,570$ (Ashkenazi), 1,711 (North African), 528 (Middle Eastern), 136 (Sephardi) and 78 (Yemenite) same ancestry pairs; $n = 61,048$ different ancestry pairs. The lower and upper limits of the boxes represent the 25% and 75% percentiles, respectively, and the top and bottom whiskers represent the 5% and 95% percentiles, respectively. **e**, Box plots of Bray–Curtis dissimilarities across pairs of 946 individuals (including admixed individuals), organized according to shared ancestry fraction (the fraction of grandparents of the same ancestry), for pairs with 0% ($n = 167,618$), 25% ($n = 33,119$), 50% ($n = 100,163$), 75% ($n = 34,187$) and 100% ($n = 111,898$) shared ancestry fractions. The lower and upper limits of the boxes represent the 25% and 75% percentiles, respectively, and the top and bottom whiskers represent the 5% and 95% percentiles, respectively. The figure demonstrates that microbiome similarity is not associated with ancestral similarity ($P = 0.73$; Mantel test).



b

Phenotype	Microbiome association index		Genetic heritability (literature)
	Israeli cohort	LLD cohort	
HDL	35.9%***	27.9%***	23.9%–48%
Lactose cons.	35.5%***	N/A	N/A
Waist circ.	28.8%***	26%***	15%–24%
Hip circ.	27.1%***	28%***	10.6%–27%
Glycaemic status	24.5%***	N/A	N/A
BMI	24.5%***	27.8%***	14%–32%
WHR	23.9%***	6.9%*	12%–14%
Fasting glucose	21.9%***	8%**	9%–33%
HbA1c%	16.1%*	8.4%	21%–32%
Creatinine	12.3%*	6.7%	19%–25%
Height	3.2%	25.9%***	33%–68%
Total cholesterol	0%	13.5%	14%–53%

indicate a greater confidence in the estimation. **b**, b^2 estimates from the analysis of 715 individuals with measured genotyped and gut microbiomes from the Israeli cohort (left column) and of 836 individuals from the LLD cohort (middle column) are comparable to previous genetic heritability estimates^{27–34} (right column). *FDR < 0.05, **FDR < 0.01 and ***FDR < 0.001. Cons., consumption, circ., circumference. **c**, b^2 estimates



***FDR < 0.001. Cons., consumption, circ., circumference. **c**, b^2 estimates of several human phenotypes and their 95% confidence intervals, evaluated using 715 individuals. *FDR < 0.05, **FDR < 0.01 and ***FDR < 0.001. **d**, Phenotype prediction accuracy for 715 individuals, evaluated using a LMM under different sets of predictive features (measured using coefficient of determination (R^2)), using four different models for each phenotype: (i) 'Basic', age, gender and diet features; (ii) 'Basic + microbiome', basic features and relative abundances of bacterial genes; (iii) 'Basic + genetics', basic features and host genotypes; and (iv) 'Basic + genetics + microbiome': basic features, relative abundances of bacterial genes and host genotypes. **e**, The additive contribution of microbiome and genetics to prediction performance evaluated using a LMM across 715 individuals, over a model that includes only basic features. The joint contribution of microbiome and genetics is similar to the sum of the individual contributions, suggesting these are independent contributions.

Basic: Age + gender + calories

- Basic
- Basic + microbiome
- Basic + genetics
- Basic + genetics + microbiome
- Microbiome
- Genetics
- Genetics + microbiome

Box 1 | Ten areas of microbiome inquiry that should be pursued

- Understanding microbiome characteristics in relation to families: which features are inherited and which are not?*
- Understanding secular trends in microbiome composition: which taxonomic groups have been lost or gained?[†]
- For diseases that have changed markedly in incidence in recent decades, do changes in the microbiome have a role? Notable examples include childhood-onset asthma, food allergies, type 1 diabetes, obesity, inflammatory bowel disease and autism.*[‡]
- Do particular signatures of the metagenome predict risks for specific human cancers and other diseases that are associated with ageing? Can these signatures be pursued to better understand oncogenesis? (Work on *Helicobacter pylori* provides a clear example of this.)*
- How do antibiotics perturb the microbiome, both in the short-term and long-term? Does the route of administration matter?*
- How does the microbiome affect the pharmacology of medications? Can we ‘micro-type’ people to improve pharmacokinetics and/or reduce toxicity? Can we manipulate the microbiome to improve pharmacokinetic stability?*[‡]
- Can we harness knowledge of microbiomes to improve diagnostics for disease status and susceptibility?*
- Can we harness the close mechanistic interactions between the microbiome and the host to provide hints for the development of new drugs?[‡]
- Specifically, can we harness the microbiome to develop new narrow-spectrum antibiotics?[‡]
- Can we use knowledge of the microbiota to develop true probiotics (and prebiotics)?*[‡]

*Areas currently under investigation. [†]Proposed areas for investigation.

ARTICLE

OPEN

doi:10.1038/nature24621

A communal catalogue reveals Earth's multiscale microbial diversity

Luke R. Thompson^{1,2,3}, Jon G. Sanders¹, Daniel McDonald¹, Amnon Amir¹, Joshua Ladau⁴, Kenneth J. Locey⁵, Robert J. Prill⁶, Anupriya Tripathi^{1,7,8}, Sean M. Gibbons^{9,10}, Gail Ackermann¹, Jose A. Navas-Molina^{1,11}, Stefan Janssen¹, Evguenia Kopylova¹, Yoshiaki Vázquez-Baeza^{1,11}, Antonio González¹, James T. Morton^{1,11}, Siavash Mirarab¹², Zhenjiang Zech Xu¹, Lingjing Jiang^{1,13}, Mohamed F. Haroon¹⁴, Jad Kanbar¹, Qiyun Zhu¹, Se Jin Song¹, Tomasz Kosciolek¹, Nicholas A. Bokulich¹⁵, Joshua Lefler¹, Colin J. Brislawn¹⁶, Gregory Humphrey¹, Sarah M. Owens¹⁷, Jarrad Hampton-Marcell^{17,18}, Donna Berg-Lyons¹⁹, Valerie McKenzie²⁰, Noah Fierer^{20,21}, Jed A. Fuhrman²², Aaron Clauset^{19,23}, Rick L. Stevens^{24,25}, Ashley Shade^{26,27,28}, Katherine S. Pollard⁴, Kelly D. Goodwin³, Janet K. Jansson¹⁶, Jack A. Gilbert^{17,29}, Rob Knight^{1,11,30} & The Earth Microbiome Project Consortium*

Our growing awareness of the microbial world's importance and diversity contrasts starkly with our limited understanding of its fundamental structure. Despite recent advances in DNA sequencing, a lack of standardized protocols and common analytical frameworks impedes comparisons among studies, hindering the development of global inferences about microbial life on Earth. Here we present a meta-analysis of microbial community samples collected by hundreds of researchers for the Earth Microbiome Project. Coordinated protocols and new analytical methods, particularly the use of exact sequences instead of clustered operational taxonomic units, enable bacterial and archaeal ribosomal RNA gene sequences to be followed across multiple studies and allow us to explore patterns of diversity at an unprecedented scale. The result is both a reference database giving global context to DNA sequence data and a framework for incorporating data from future studies, fostering increasingly complete characterization of Earth's microbial diversity.



BY THE NUMBERS

27,751

samples

7 continents

43 countries

2,212,796,183

total DNA sequences

307,572

unique DNA sequences
(approx. species)

50+

peer-reviewed
publications

500+

scientists



92 environmental
features

2 – 12 pH range

66 animal host
species

(stomach acid to household
ammonia)

78.9 °N –
78.2 °S

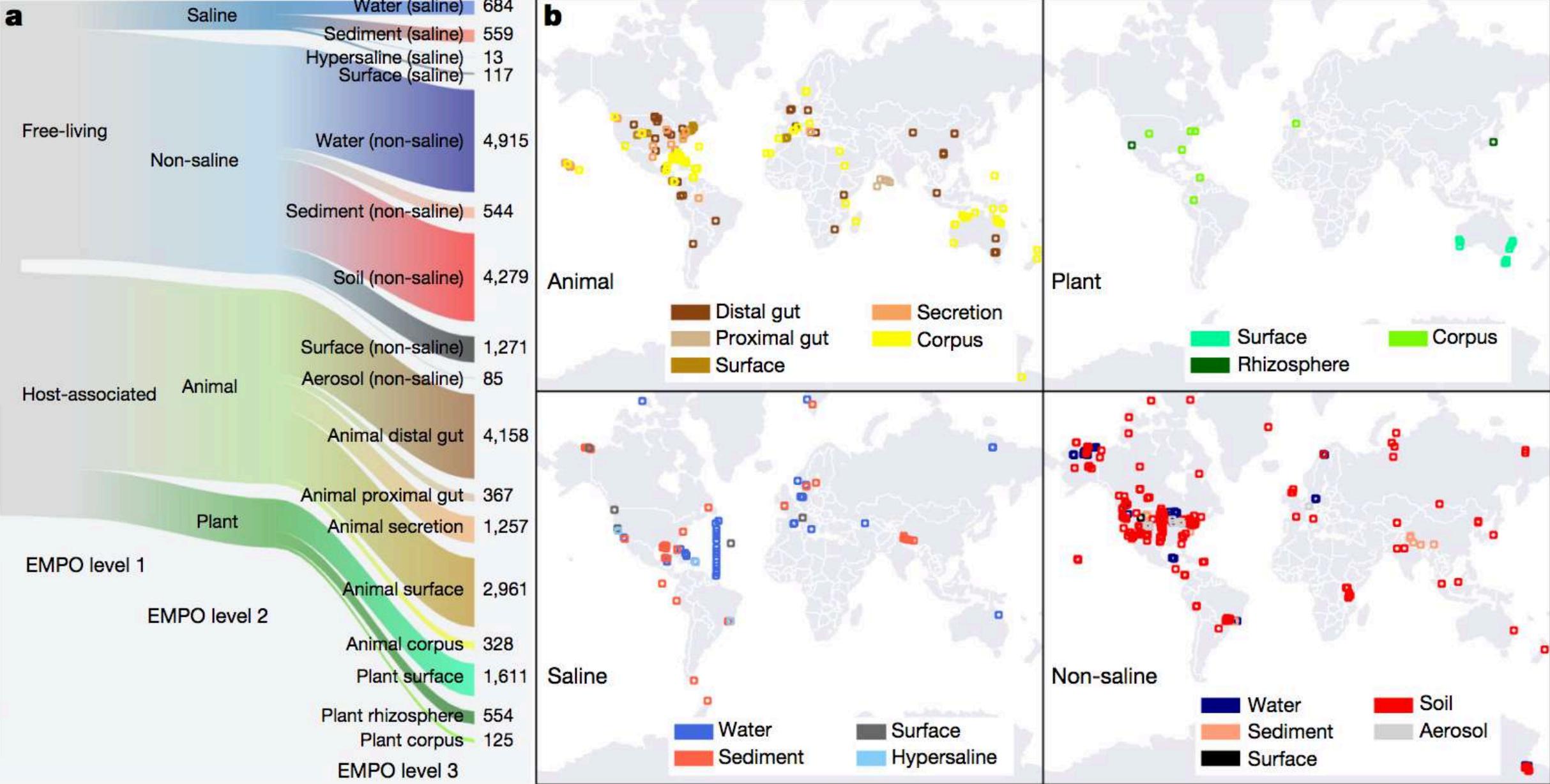
latitude range (Arctic
Circle to Antarctica)

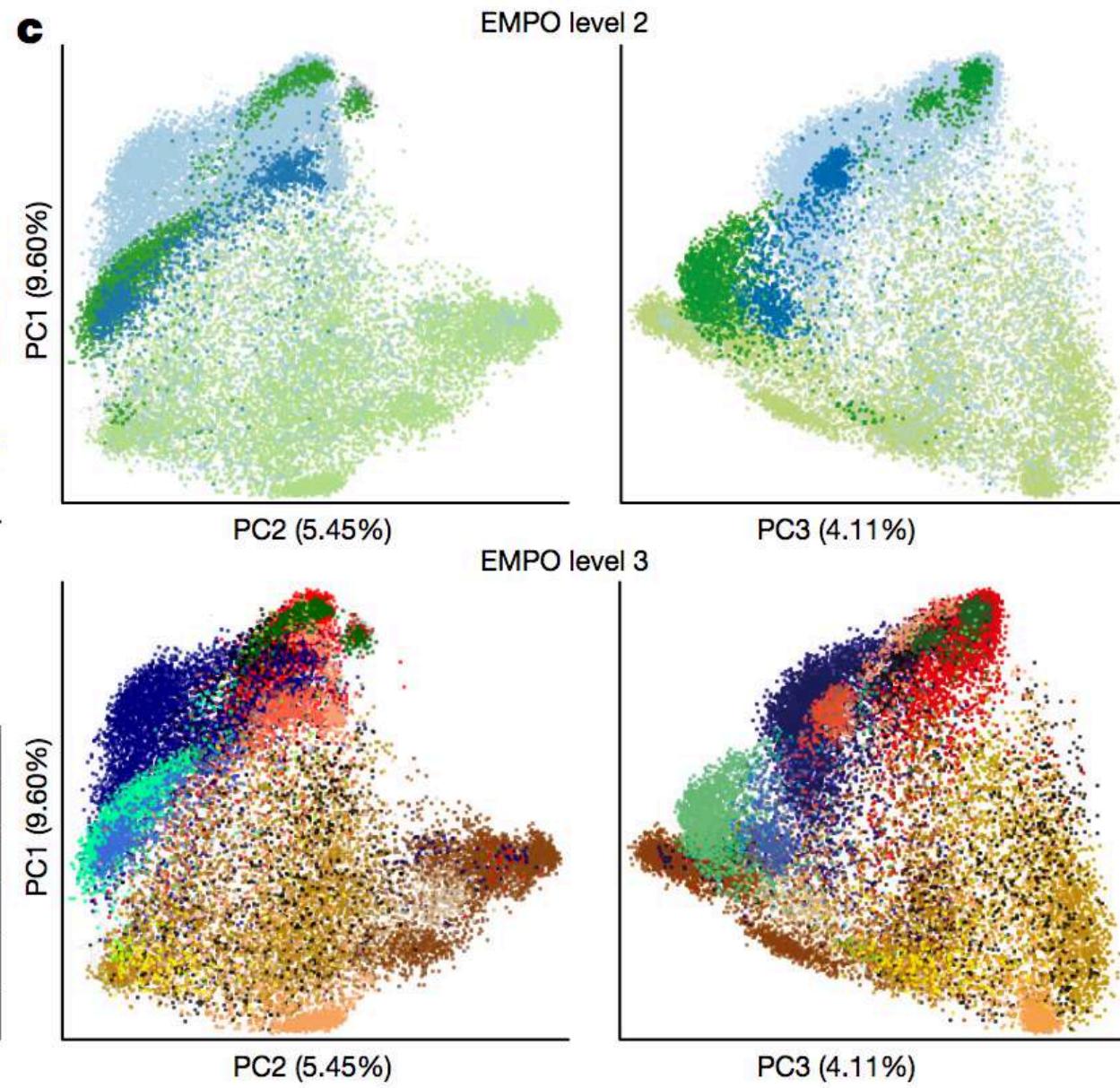
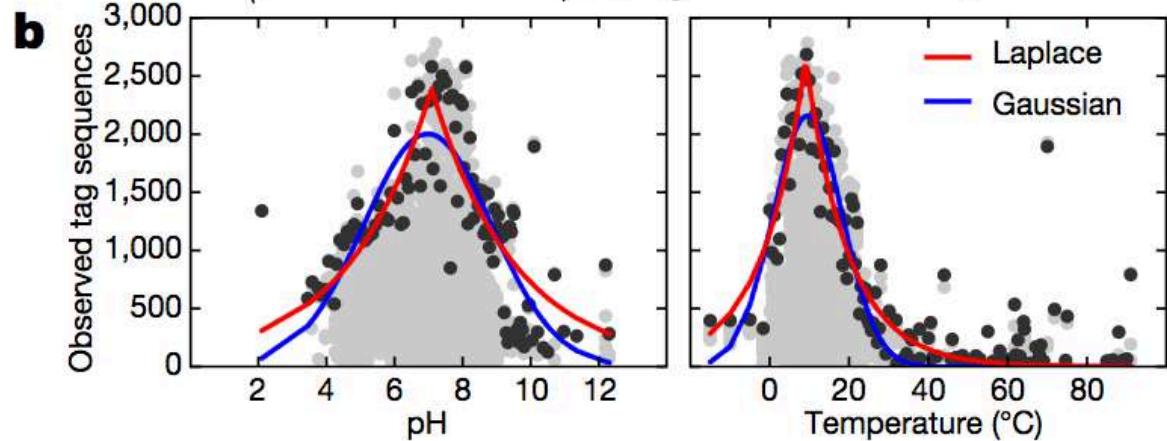
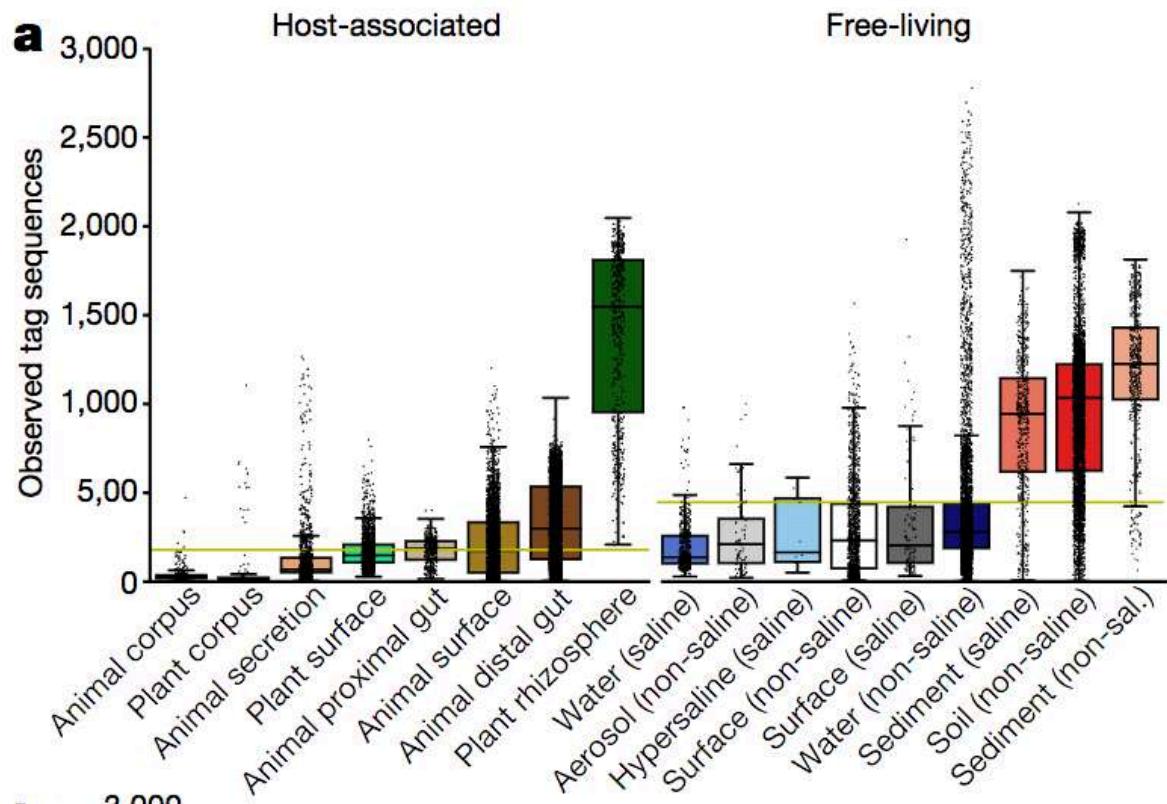
1

reference database of
bacteria that reside on
Planet Earth

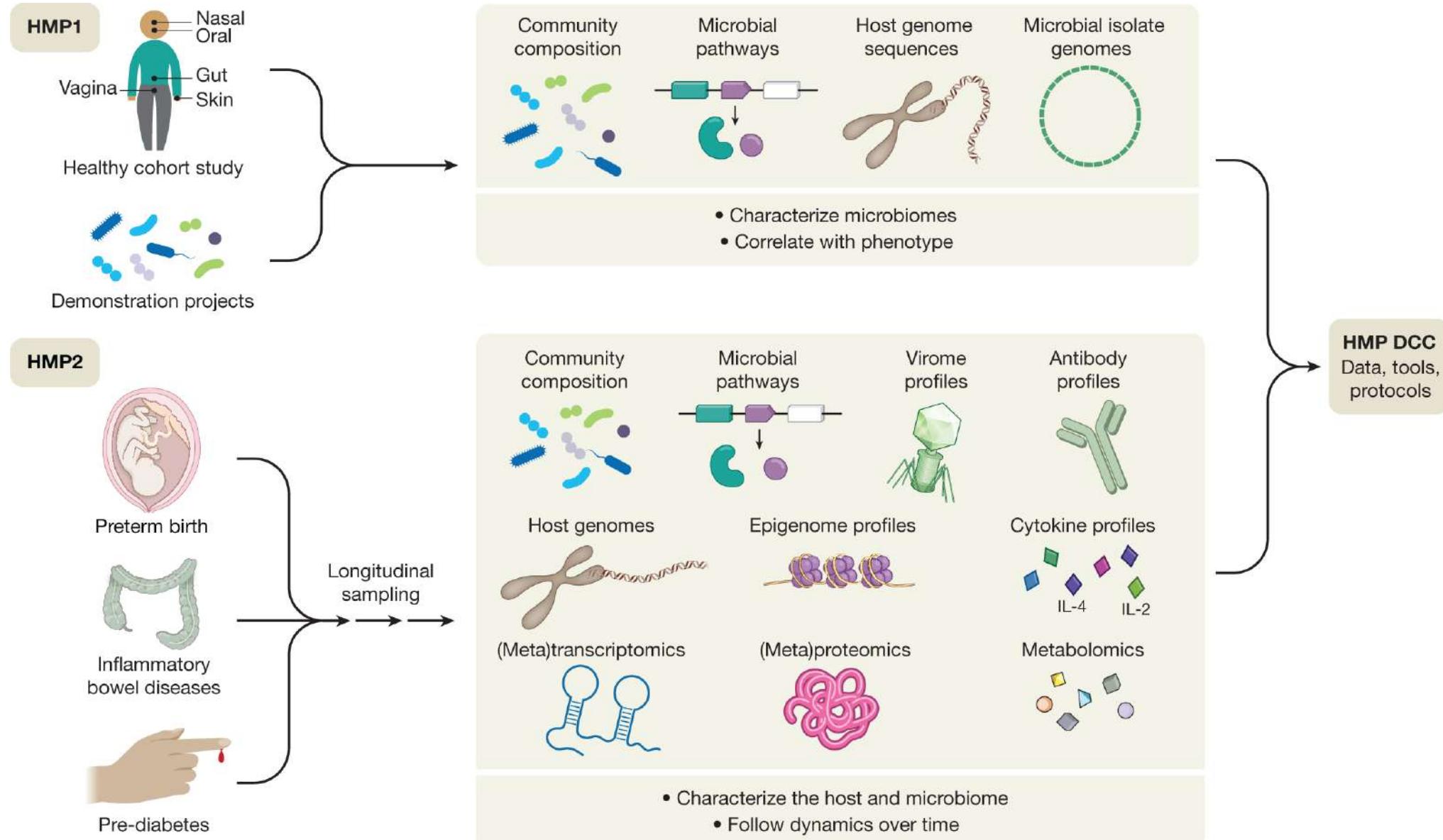
Credit: UC San Diego Center for Microbiome Innovation

Reference: Thompson et al., *Nature*, 2017. doi:10.1038/nature24621





“a paradigm for future multi-omic studies of the human microbiome”



New challenges

- So much data
- Technology advancement
- **Integrating different kinds of data (multi-omic)**
- High performance
- Reproducibility crisis
- Bioinformaticians as a profession
- Only biology has a specific term to refer to the use of computers in this discipline ('bioinformatics')
- Proper integration into academic curriculums