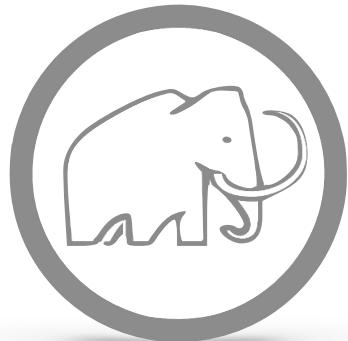


# Interactive Microbiome Analysis using **DIAMOND+MEGAN**



Daniel H. Huson  
Banu Cetinkaya



Tutorial web page

# Outline

- Introduction to microbiome analysis
- Step 0: Software setup
- Step 1: DIAMOND alignment against protein database
- Step 2: MEGANization of reads and alignments
- Step 3: MEGAN interactive analysis

# Outline

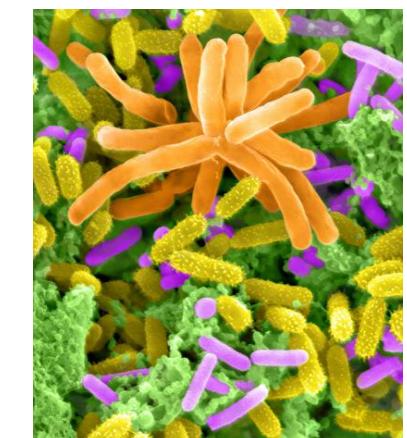
- Introduction to microbiome analysis
- Step 0: Installation
- Step 1: DIAMOND alignment against protein database
- Step 2: MEGANization of reads and alignments
- Step 3: MEGAN interactive analysis

# Microbiome

- Traditionally, microbes are studied in pure culture
- **Genome:**
  - Entire DNA sequence of a single organism
- *But:* most microbes don't live in isolation and many can't be cultured
- **Microbiome:**
  - Collection of microbes in a specific theatre of activity
- **Metagenome:**
  - Entire DNA sequence of a microbiome



[www.innovations-report.de](http://www.innovations-report.de)



[www.physorg.com](http://www.physorg.com)

# Sources of studied microbiomes

- Soil samples
- Water samples
- Seabed samples
- Air samples
- Ancient bones
- Host-associated samples
- Human microbiome
- ....



soils.usda.gov



<http://outdoors.webshots.com>



CSIRO  
www.scienceimage.csiro.au



www.lanl.gov



- First NGS technique 454 released
- Intended for genome sequencing...

nature  
Vol 437 | 15 September 2005 doi:10.1038/nature03959

## ARTICLES

### Genome sequencing in microfabricated high-density picolitre reactors

Marcel Margulies<sup>1,\*</sup>, Michael Egholm<sup>1,\*</sup>, William E. Altman<sup>1</sup>, Said Attiya<sup>1</sup>, Joel S. Bader<sup>1</sup>, Lisa A. Bemben<sup>1</sup>, Jan Berka<sup>1</sup>, Michael S. Braverman<sup>1</sup>, Yi-Ju Chen<sup>1</sup>, Zhoutao Chen<sup>1</sup>, Scott B. Dewell<sup>1</sup>, Lei Du<sup>1</sup>, Joseph M. Fierro<sup>1</sup>, Xavier V. Gomes<sup>1</sup>, Brian C. Godwin<sup>1</sup>, Wen He<sup>1</sup>, Scott Helgesen<sup>1</sup>, Chun He Ho<sup>1</sup>, Gerard P. Iرزک<sup>1</sup>, Szilveszter C. Jando<sup>1</sup>, Maria L. I. Alequer<sup>1</sup>, Thomas P. Jarvie<sup>1</sup>, Kshama B. Jirage<sup>1</sup>, Jong-Bum Kim<sup>1</sup>, James R. Knight<sup>1</sup>, Janna R. Lanza<sup>1</sup>, John H. Leamon<sup>1</sup>, Steven M. Lefkowitz<sup>1</sup>, Ming Lei<sup>1</sup>, Jing Li<sup>1</sup>, Kenton L. Lohman<sup>1</sup>, Hong Lu<sup>1</sup>, Vinod B. Makrilia<sup>1</sup>, Keith E. McDade<sup>1</sup>, Michael P. McKenna<sup>1</sup>, Eugene W. Myers<sup>2</sup>, Elizabeth Nickerson<sup>1</sup>, John R. Nobile<sup>1</sup>, Ramona Plant<sup>1</sup>, Bernard P. Puc<sup>1</sup>, Michael T. Ronan<sup>1</sup>, George T. Roth<sup>1</sup>, Gary J. Sarkis<sup>1</sup>, Jan Fredrik Simons<sup>1</sup>, John W. Simpson<sup>1</sup>, Maithreyan Srinivasan<sup>1</sup>, Karrie R. Tartar<sup>1</sup>, Alexander Tomasz<sup>1</sup>, Kari A. Vogt<sup>1</sup>, Greg A. Volkmer<sup>1</sup>, Shally H. Wang<sup>1</sup>, Yong Wang<sup>1</sup>, Michael P. Weiner<sup>4</sup>, Penruoane Yu<sup>1</sup>, Richard F. Bewley<sup>1</sup> & Jonathan M. Rothberg<sup>1</sup>



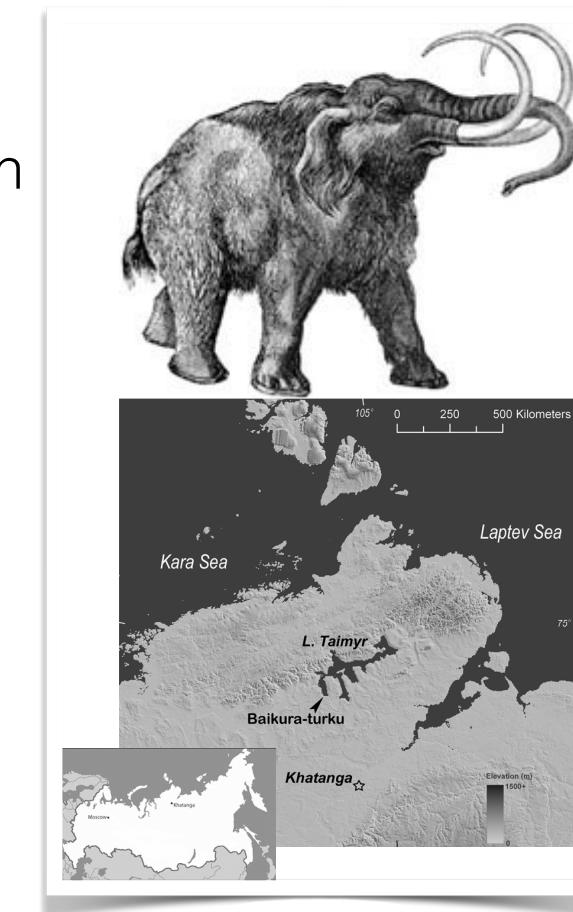
★Use NGS to sequence ancient DNA?

★Use NGS to sequence metagenomic DNA?

NGS = next generation sequencing

# Mammoth DNA & metagenome (2006)

- DNA collected from permafrost mammoth (28,000 years old)
- DNA extracted from 1g bone
- DNA sheared to 500-700 bp
- Sequenced using 454
- ~302,000 reads, length ~95 bp



- ★ Can use NGS for ancient DNA
- ★ First NGS metagenomics paper

## REPORTS

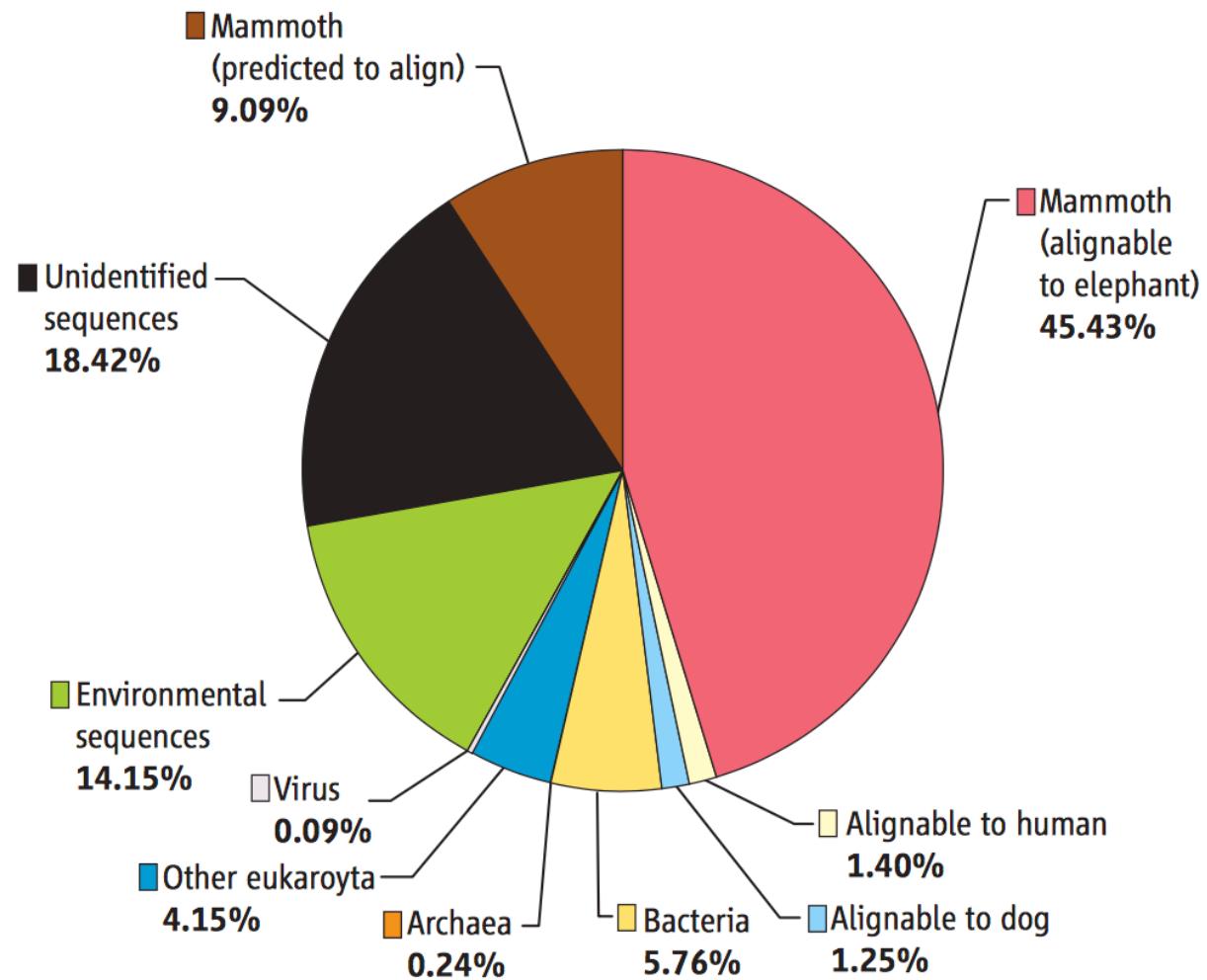
### Metagenomics to Paleogenomics: Large-Scale Sequencing of Mammoth DNA

Hendrik N. Poinar,<sup>1,2,3\*</sup> Carsten Schwarz,<sup>1,2</sup> Ji Qi,<sup>4</sup> Beth Shapiro,<sup>5</sup> Ross D. E. MacPhee,<sup>6</sup> Bernard Buigues,<sup>7</sup> Alexei Tikhonov,<sup>8</sup> Daniel H. Huson,<sup>9</sup> Lynn P. Tomsho,<sup>4</sup> Alexander Auch,<sup>9</sup> Markus Rappaport,<sup>10</sup> Webb Miller,<sup>4</sup> Stephan C. Schuster<sup>4\*</sup>

Science, 2006

# Mammoth bone metagenome (2006)

**Fig. 1.** Characterization of the mammoth metagenomic library, including percentage of read distributions to various taxa. Host organism prediction based on BLASTZ comparison against GenBank and environmental sequences database.

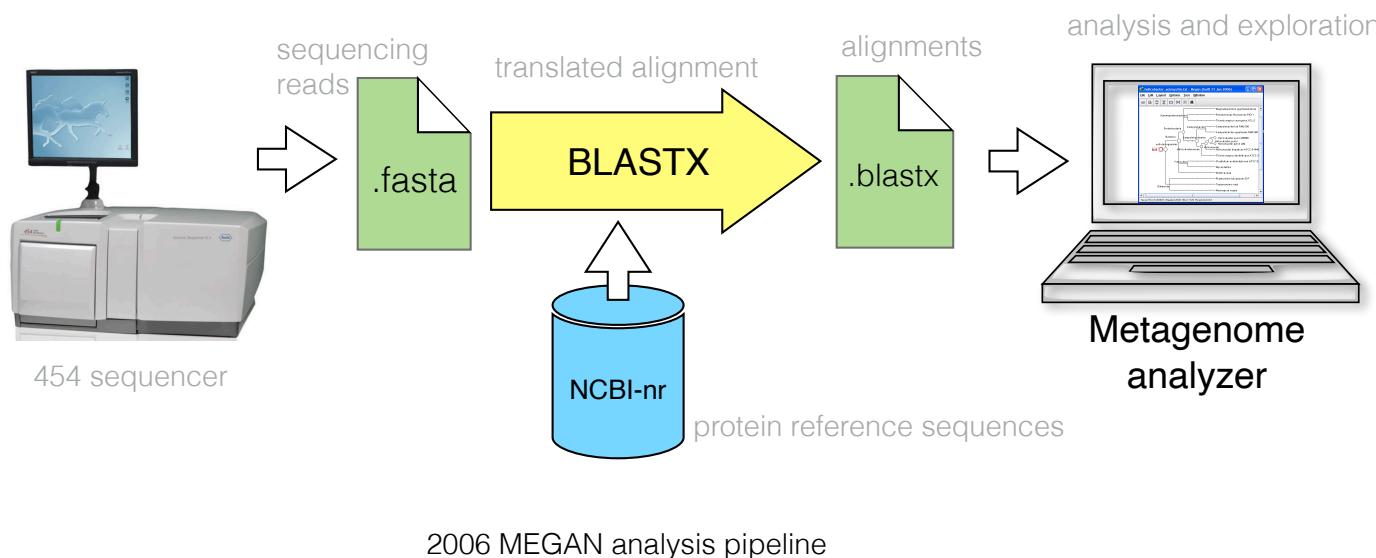


Poinar et al, Science 2006

# How to analyze metagenomic reads? (2006)

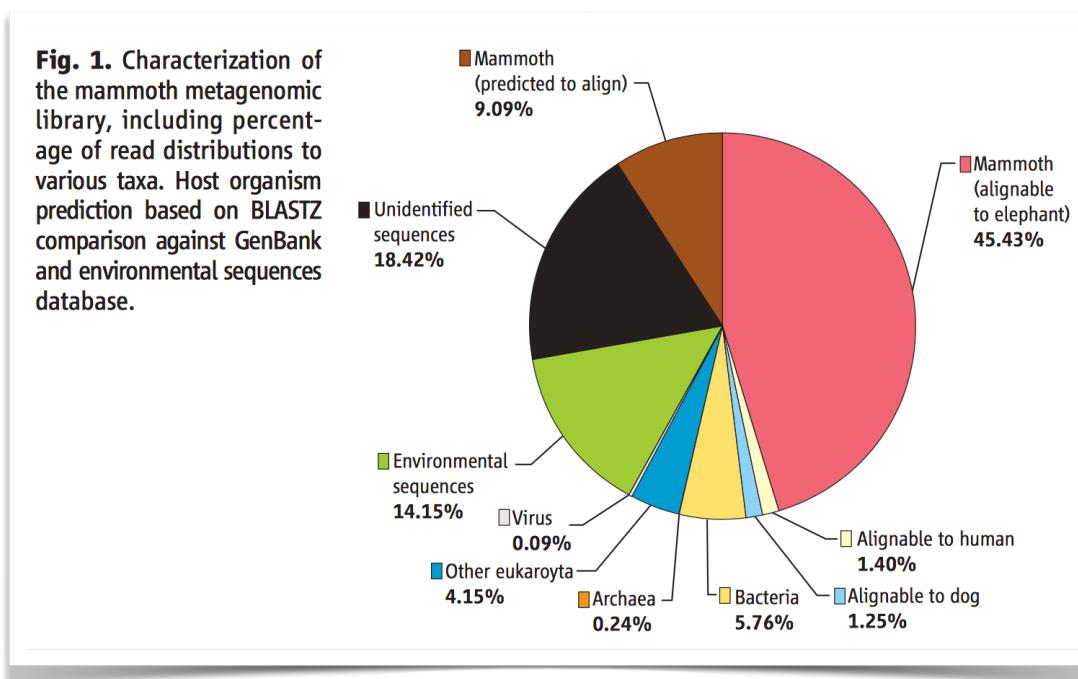
Basic idea (with Stephan Schuster at Penn State):

- BLASTX non-host reads against NCBI-nr
- Assign reads to NCBI taxonomy using naive LCA (lowest common ancestor) approach
- Develop GUI to explore assignments and alignments

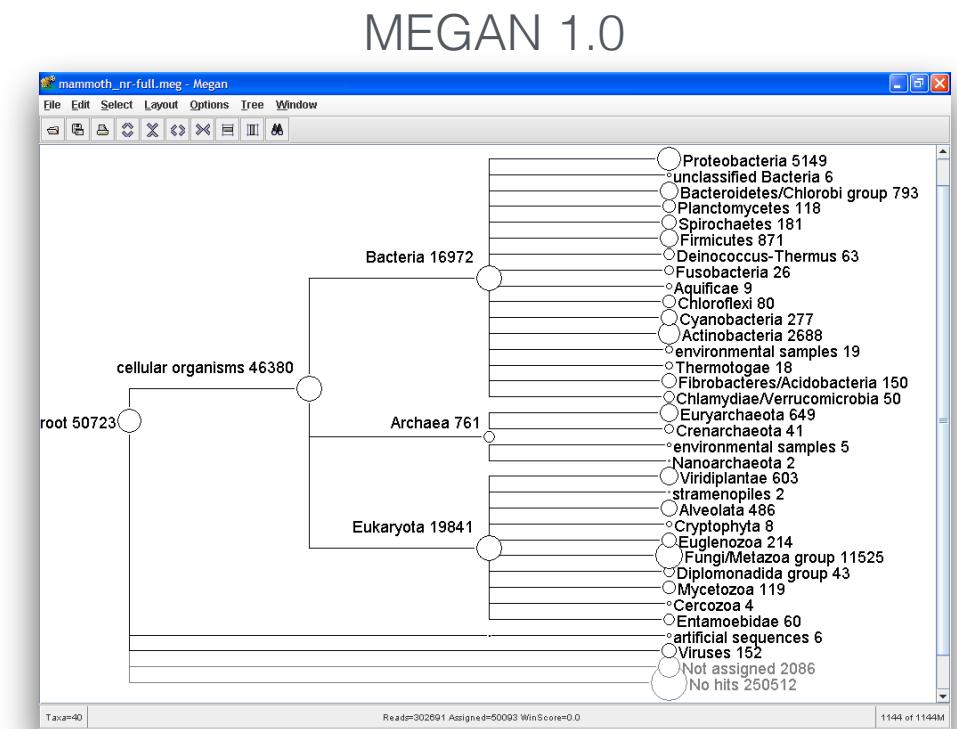


# How to analyze metagenomic reads? (2006)

- MEGAN (MEtagenome ANalyzer 1.0)



Poinar et al, Science 2006



H. et al, Genome Research, 2007

# Computational bottleneck (2006)

- Compare all reads against the NCBI-nr protein database
- Year 2006:
  - 300,000 reads of length ~100bp
  - NCBI-nr: 3 million entries, ~1 billion letters

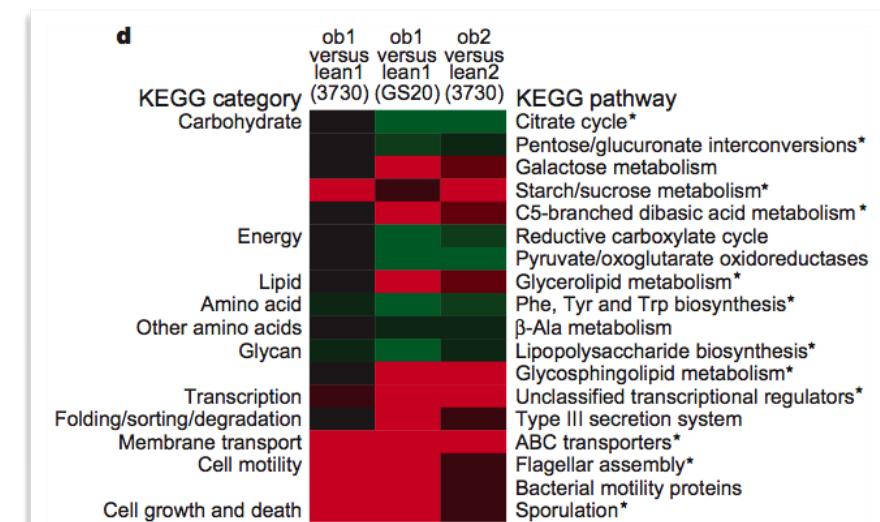
★ BLASTX took a couple of weeks on a small cluster

(NCBI-nr today: ~ 550 million entries)

# Obesity-associated gut microbiome

Turnbaugh *et al* (2006):

- Caecal microbial DNA of ob/ob, ob/+, +/+ mice
- Sanger sequencing:
  - 39.5 Mb
  - read length 750 bp
- 454 sequencing:
  - 160 Mb
  - read length 93 bp
- Change in relative abundance of Bacteroidetes and Firmicutes
- Change in functional capacity (toward energy harvesting)



# Large-scale human gut analysis

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

nature

## MetaHIT 2010

## ARTICLES

### A human gut microbial gene catalogue established by metagenomic sequencing

Junjie Qin<sup>1\*</sup>, Ruiqiang Li<sup>1\*</sup>, Jeroen Raes<sup>2,3</sup>, Manimozhiyan Arumugam<sup>2</sup>, Kristoffer Solvsten Burgdorf<sup>4</sup>, Chaysavanh Manichanh<sup>5</sup>, Trine Nielsen<sup>4</sup>, Nicolas Pons<sup>6</sup>, Florence Levenez<sup>6</sup>, Takuji Yamada<sup>2</sup>, Daniel R. Mende<sup>2</sup>, Junhua Li<sup>1,7</sup>, Junming Xu<sup>1</sup>, Shaochuan Li<sup>1</sup>, Dongfang Li<sup>1,8</sup>, Jianjun Cao<sup>1</sup>, Bo Wang<sup>1</sup>, Huiqing Liang<sup>1</sup>, Huisong Zheng<sup>1</sup>, Yinlong Xie<sup>1,7</sup>, Julien Tap<sup>6</sup>, Patricia Lepage<sup>6</sup>, Marcelo Bertalan<sup>9</sup>, Jean-Michel Batto<sup>6</sup>, Torben Hansen<sup>4</sup>, Denis Le Paslier<sup>10</sup>, Allan Linneberg<sup>11</sup>, H. Bjørn Nielsen<sup>9</sup>, Eric Pelletier<sup>10</sup>, Pierre Renault<sup>6</sup>, Thomas Sicheritz-Ponten<sup>9</sup>, Keith Turner<sup>12</sup>, Hongmei Zhu<sup>1</sup>, Chang Yu<sup>1</sup>, Shengting Li<sup>1</sup>, Min Jian<sup>1</sup>, Yan Zhou<sup>1</sup>, Yingrui Li<sup>1</sup>, Xiuqing Zhang<sup>1</sup>, Songgang Li<sup>1</sup>, Nan Qin<sup>1</sup>, Huanming Yang<sup>1</sup>, Jian Wang<sup>1</sup>, Søren Brunak<sup>9</sup>, Joel Doré<sup>6</sup>, Francisco Guarner<sup>5</sup>, Karsten Kristiansen<sup>13</sup>, Oluf Pedersen<sup>4,14</sup>, Julian Parkhill<sup>12</sup>, Jean Weissenbach<sup>10</sup>, MetaHIT Consortium†, Peer Bork<sup>2</sup>, S. Dusko Ehrlich<sup>6</sup> & Jun Wang<sup>1,13</sup>

To understand the impact of gut microbes on human health and well-being it is crucial to assess their genetic potential. Here we describe the Illumina-based metagenomic sequencing, assembly and characterization of 3.3 million non-redundant microbial genes, derived from 576.7 gigabases of sequence, from faecal samples of 124 European individuals. The gene set, ~150 times larger than the human gene complement, contains an overwhelming majority of the prevalent (more frequent) microbial genes of the cohort and probably includes a large proportion of the prevalent human intestinal microbial genes. The genes are largely shared among individuals of the cohort. Over 99% of the genes are bacterial, indicating that the entire cohort harbours between 1,000 and 1,150 prevalent bacterial species and each individual at least 160 such species, which are also largely shared. We define and describe the minimal gut metagenome and the minimal gut bacterial genome in terms of functions present in all individuals and most bacteria, respectively.

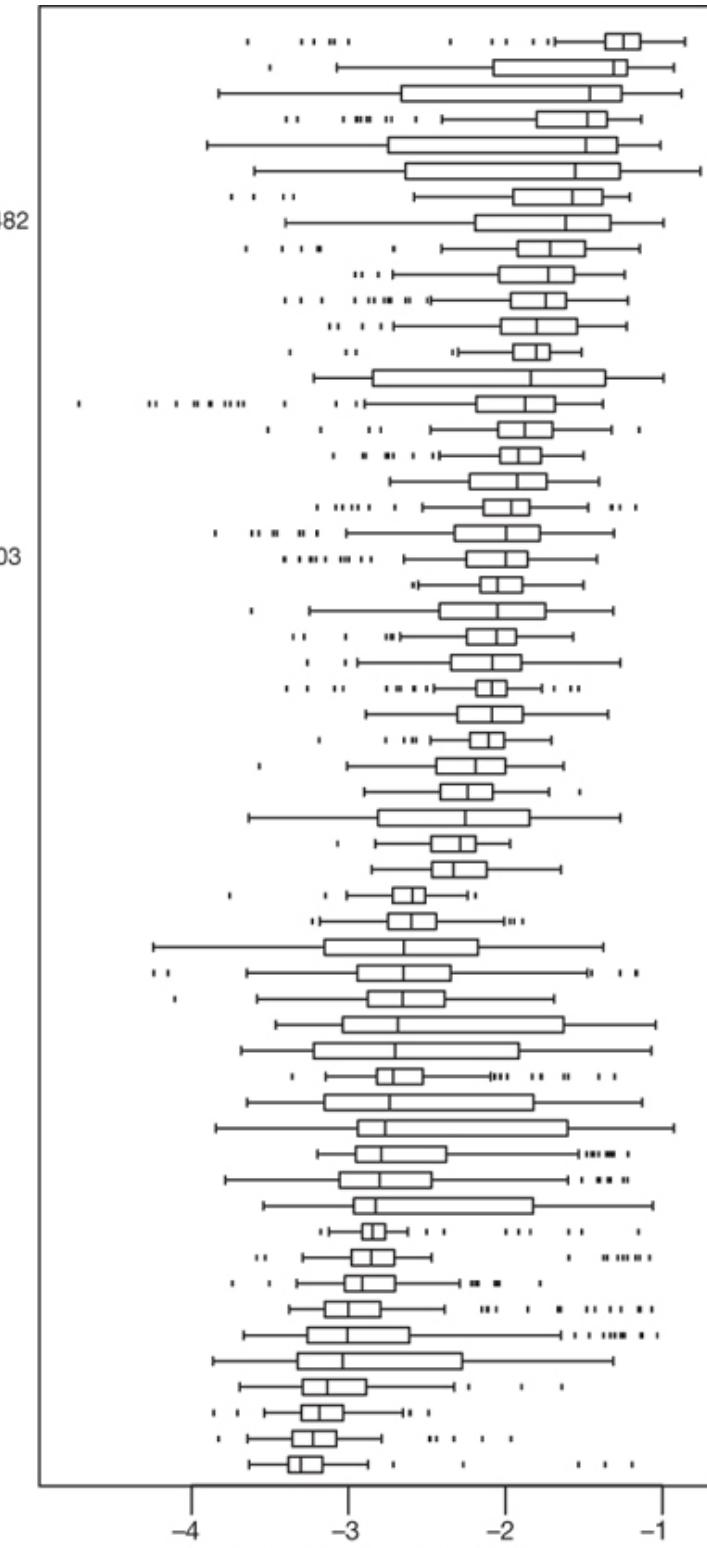
- 576Gb of sequence from 124 individuals

# Core of human gut microbiome

- 57 species present in  $\geq 90\%$  of individuals with coverage  $> 1\%$
- High variability
- Bacteroidetes and Firmicutes most abundant

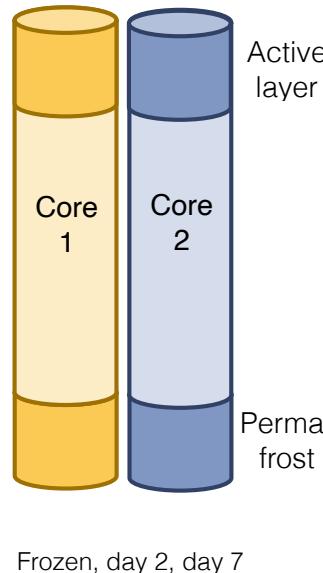
BLASTX at Super Computer Center in Barcelona, then MEGAN analysis

*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 formicigenans*  
*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 obaeum* A2–162  
*Subdoligranulum variabile*  
*Bacteroides capillosus*  
*Streptococcus thermophilus* LMD–9  
*Clostridium leptum*  
*Holdemania filiformis*  
*Bacteroides stercoris*  
*Coprococcus eutactus*  
*Clostridium* sp. M62 1  
*Bacteroides eggerthii*  
*Butyrivibrio crossotus*  
*Bacteroides finegoldii*  
*Parabacteroides johnsonii*  
*Clostridium* sp. L2–50  
*Clostridium nexile*  
*Bacteroides pectinophilus*  
*Anaerotruncus colihominis*  
*Streptococcus gnavus*  
*Bacteroides intestinalis*  
*Bacteroides fragilis* 3\_1\_12  
*Clostridium asparagiforme*  
*Enterococcus faecalis* TX0104  
*Clostridium scindens*  
*Blautia hansenii*

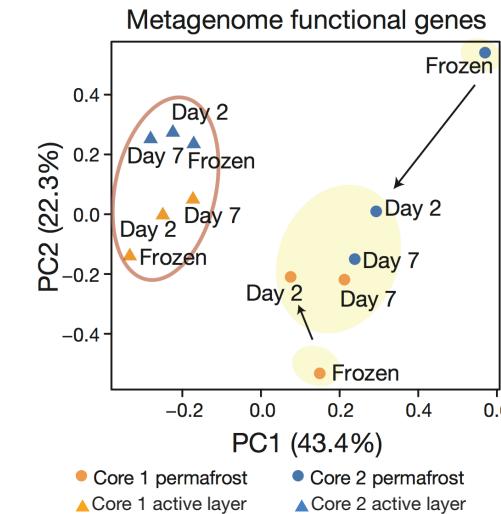


# Permafrost study (2011)

(Mackelprang *et al*, Science 2011)



Their question:  
Functional changes  
during thawing?



- Align ~250 million Illumina reads against KEGG
- 800,000 CPU hours at Super Computer Center in Berkeley



on 100 cores

# Three basic questions

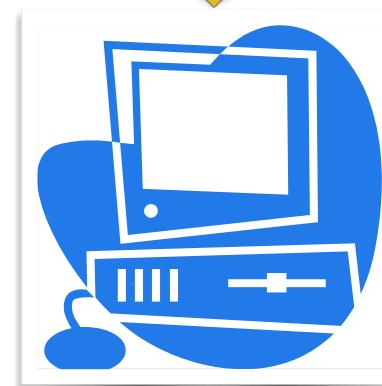


Hundreds of Samples



High-throughput  
DNA sequencing

Billions of sequences

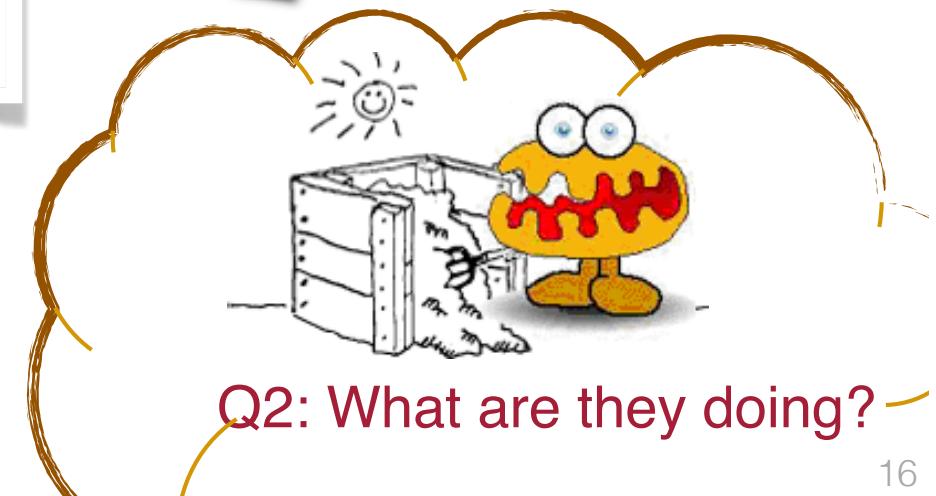


Basic computational  
analysis

Many  
CPU hours

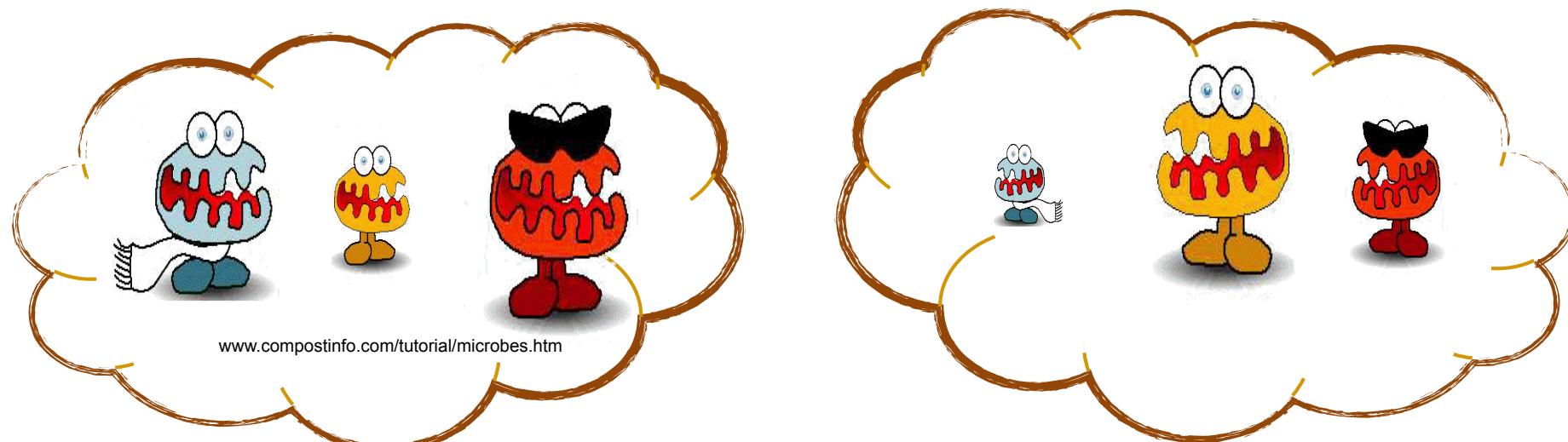


Q1: Who is out there?



Q2: What are they doing? —

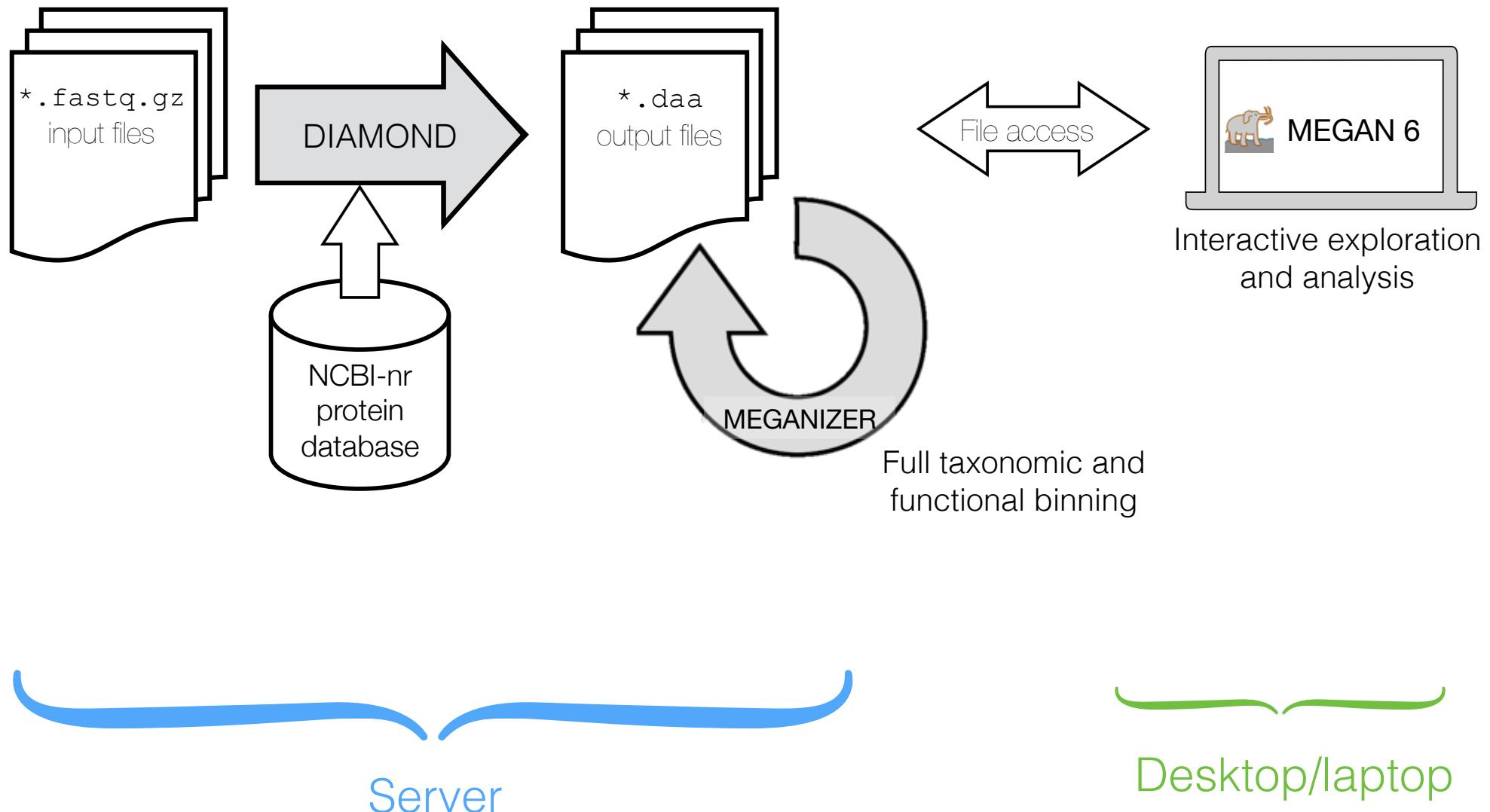
# Q3: How do they compare?



# Three-step approach

- Step 1: Align reads or assembled contigs against protein reference sequences - DIAMOND
- Step 2: Analyze alignments to assign sequences to taxonomic and functional classes - MEGANIZER
- Step 3: Interactively explore, analyze and compare samples - MEGAN

# DIAMOND+MEGAN pipeline



# Outline

- Introduction to microbiome analysis
- Step 0: Software setup
- Step 1: DIAMOND alignment against protein database
- Step 2: MEGANization of reads and alignments
- Step 3: MEGAN interactive analysis

# Step 0 - Installation

- Tutorial webpage, available from:

<https://github.com/husonlab/tutorials/wiki/Tutorial>



You should have the following files and software installed:

- DIAMOND, from:

<https://github.com/bbuchfink/diamond>

- MEGAN, from:

<https://software-ab.cs.uni-tuebingen.de/download/megan6>

- Data files, tutorial reference and mapping files, from:

<https://software-ab.cs.uni-tuebingen.de/download/megan6/tutorial/tutorial-dm.zip>

# Outline

- Introduction to microbiome analysis
- Step 0: Installation
- Step 1: DIAMOND alignment against protein database
- Step 2: MEGANization of reads and alignments
- Step 3: MEGAN interactive analysis

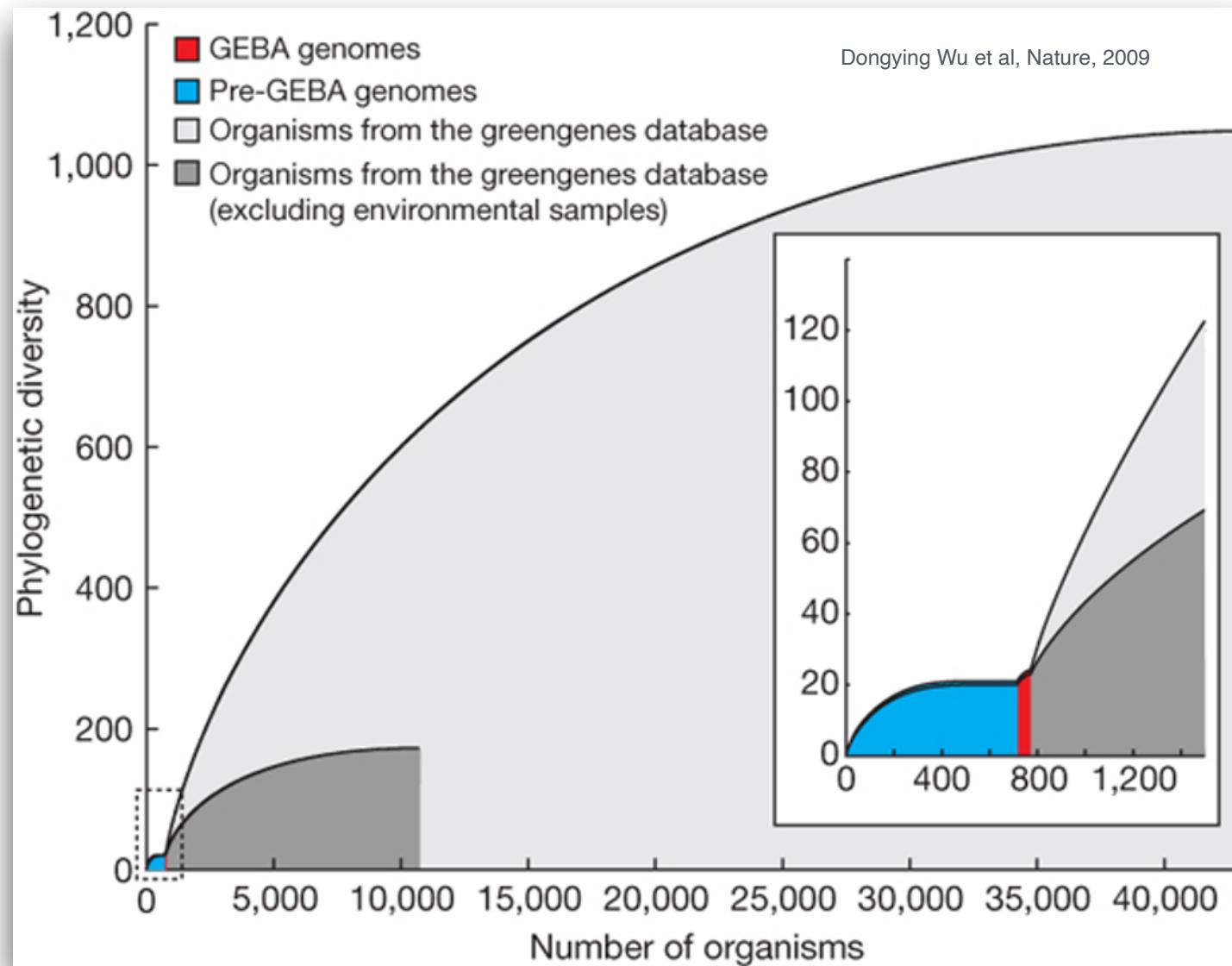
# Step I - Protein alignment

- Key idea: To perform metagenome analysis, align against a protein reference database

Why protein alignment?

- To *identify known* genomes in a sample, use DNA alignment (e.g., pathogen detection, human gut)
- To analyze *unknown* organisms, protein alignment is more suitable due to higher level of conservation

# Genome databases don't cover enough diversity



# Translated alignment

- Read:

>HISEQ:457:C5366ACXX:2:1101:5937:60460 (101 bases)  
**TTATATTAATTAGAAAACCAATTAAAAATACGAACGTTATGAAGAAGTACATTGC...**

- Translation (frame +3):

. . I   L   I   R   K   P   I   K   N   T   N   V   M   K   K   Y   I   C   ...

- Translated alignment:

>EEC52678.1 Length = 65

Score = 56 bits (135), Expect = 1e-05

Identities = 22/33 (67%), Positives = 27/33 (82%), Gaps = 0/33 (0%)

Frame = +3

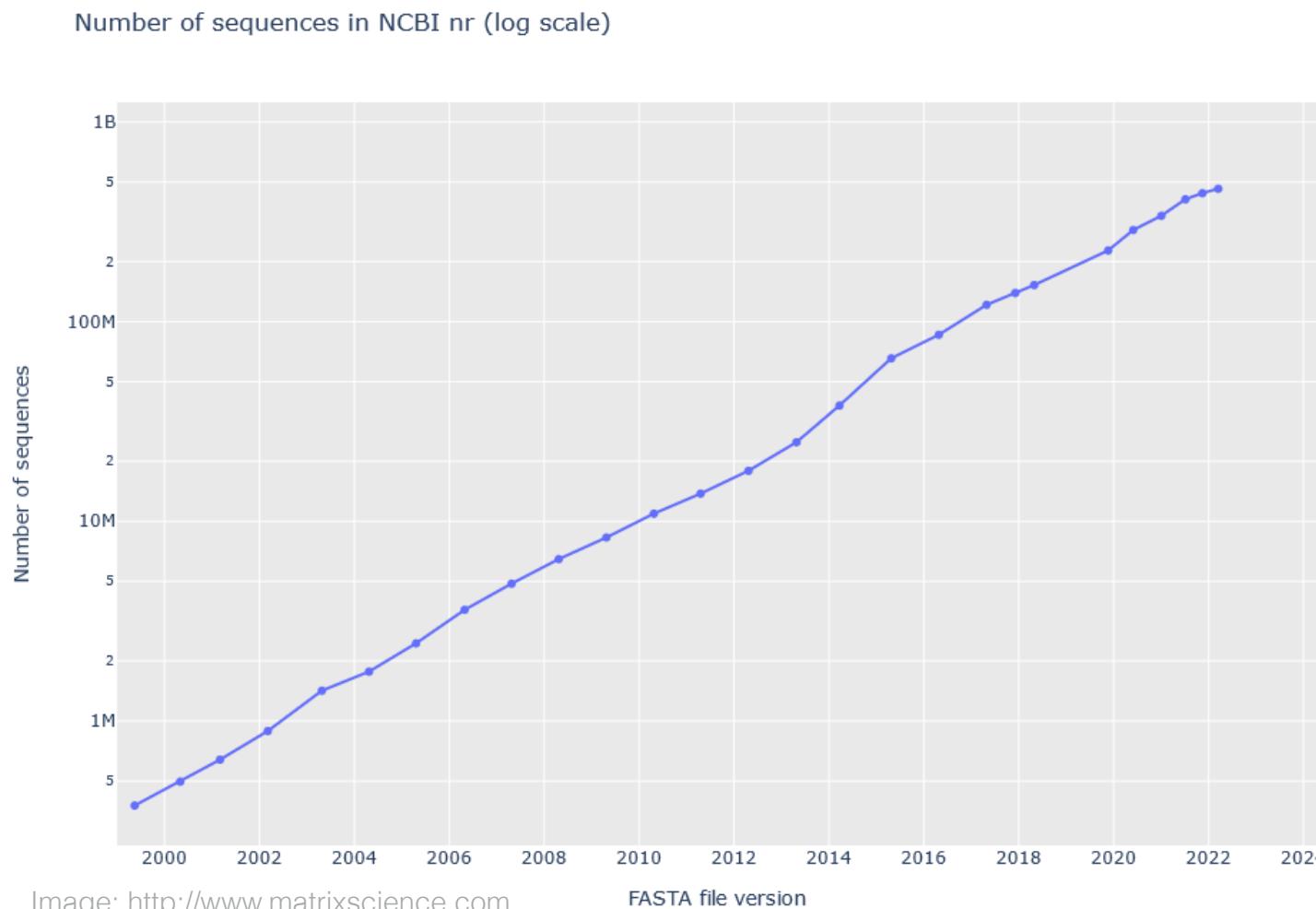
Query:            3    **ILIRKPIKNTNVMKYICTVCEYIYDPEQGDPE**    101

                  +L +K    K    VM+KYICT+CEY+YDPEQGDPE

Sbjct:            1    **MLSKKKFKQKRVMEKYICTICEYVYDPEQGDPE**    33

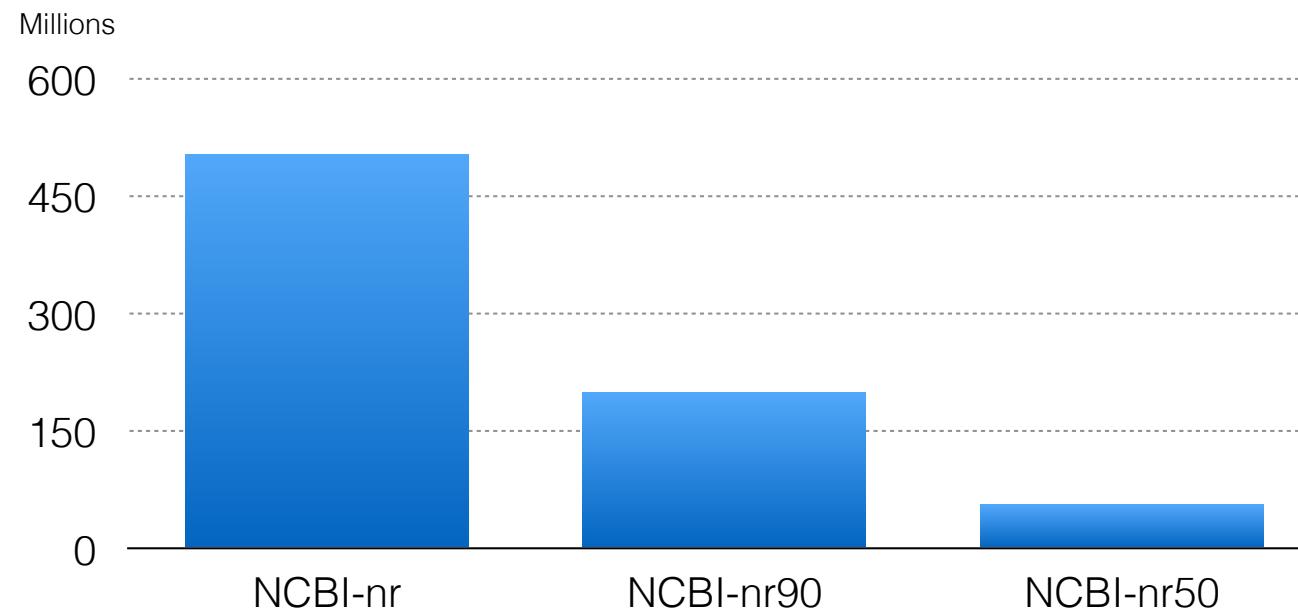
# Comprehensive database: NCBI-nr

- NCBI-nr database of non-redundant protein sequences has over 500M entries



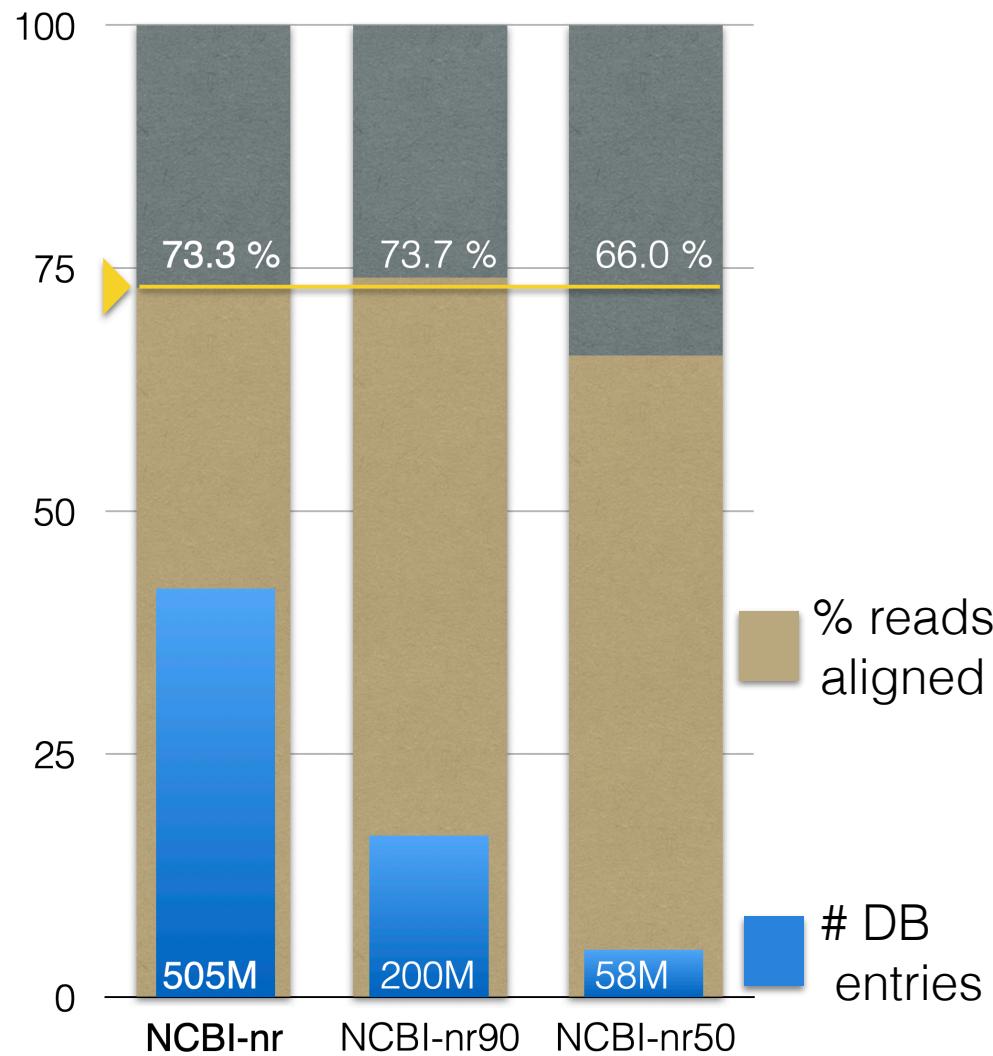
# Smaller alternative?

- Latest release of DIAMOND allows clustering of NCBI-nr (Buchfink et al, submitted)
- Similarity: 90 and 50%
- NCBI-nr90 and NCBI-nr50



# Comparison NCBI-nr vs nr50

- 201M reads aligned by DIAMOND & analyzed using MEGAN



Classification	NCBI-nr	NCBI-90	NCBI-50
NCBI			
Taxonomy	1.0	1.01	0.85
GTDB			
Taxonomy	1.0	1.03	0.88
INTERPRO	1.0	1.34	1.39
SEED	1.0	1.07	1.23
EC	1.0	1.12	1.17
EGGNOG	1.0	1.55	1.87
Speed-up	1.0	3.0	34.6

# Tutorial-nr file

- While NCBI-nr50 is small enough to be used on a high-end laptop, it is too big for this tutorial
- For the tutorial, we provide `tutorial-nr.gz`
  - This is a small subset of `nr50.gz`
  - It only contains accessions relevant for the provided short-read datasets
  - It can not be used for analysis of real data

# Step I - Protein alignment

- Will use DIAMOND
  - designed for metagenomics (Buchfink et al, 2015)
- Need to:
  - Install DIAMOND
  - Download reference sequence
  - Run DIAMOND to build index
  - Run DIAMOND on fastq.gz files

# Step I - DIAMOND index

- Build a DIAMOND index:

```
diamond makedb --in tutorial-nr.gz -d tutorial-nr
```

- Note: Using `tutorial-nr.gz`, due to time restrictions

# Step I - Run DIAMOND

- Run DIAMOND on one input FASTQ file:

```
diamond blastx -d tutorial-nr \
-q data/Alice00-1mio.fq.gz \
-o out/Alice00-1mio.daa \
-f 100 --masking 0
```

- Run DIAMOND on *all* input files in the directory:

```
for file in data/*.fq.gz
do
ofile="out/${basename "${file%.*}}.daa"
diamond blastx --db tutorial-nr \
-q $file -o $ofile -f 100 --masking 0
done
```

# Step I - Run DIAMOND

- For full size datasets, DIAMOND alignment (and subsequent meganization) is run on a server
- The 12 small datasets against `tutorial-nr.gz` will take less than 20 minutes on a modern laptop
- If you failed to run DIAMOND on the data, you can download the resulting files here:

[https://software-ab.cs.uni-tuebingen.de/download/  
megan6/tutorial/diamond-out.zip](https://software-ab.cs.uni-tuebingen.de/download/megan6/tutorial/diamond-out.zip)

# Outline

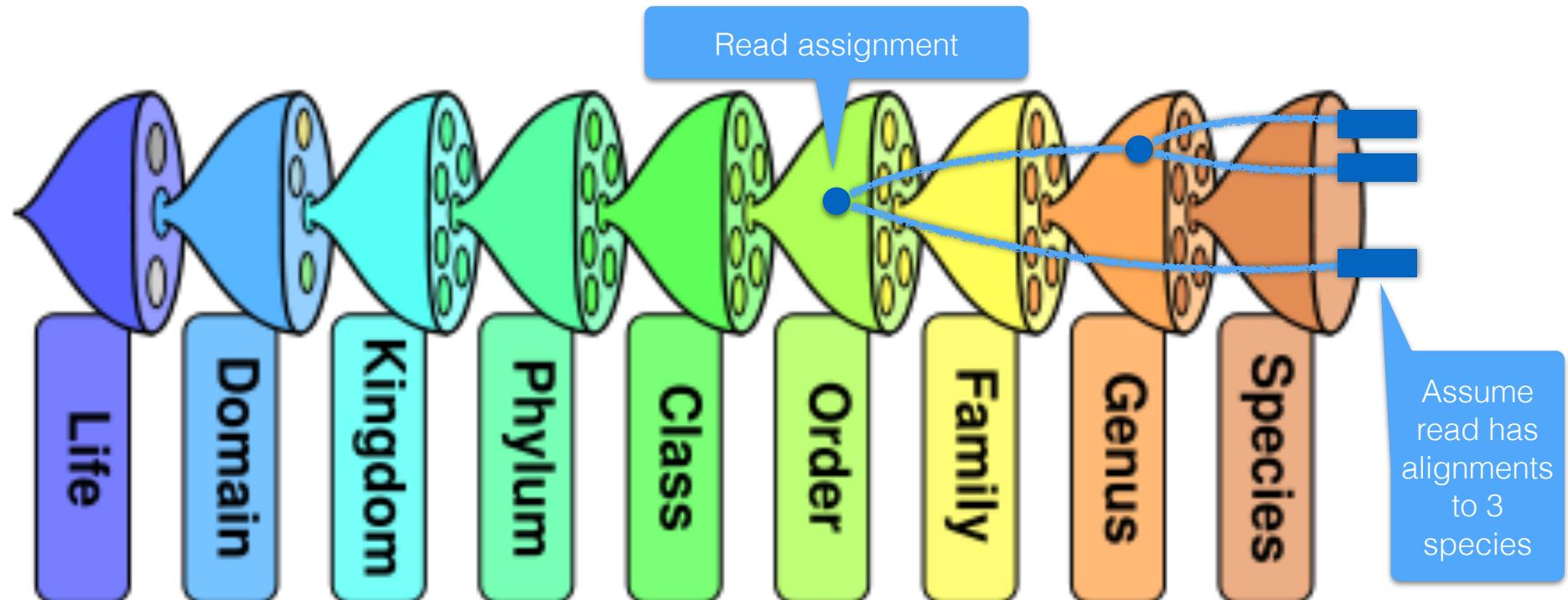
- Introduction to microbiome analysis
- Step 0: Installation
- Step 1: DIAMOND alignment against protein database
- Step 2: MEGANization of reads and alignments
- Step 3: MEGAN interactive analysis

# Step 2: Meganization

- Run DIAMOND with option -f 100, so that
  - the output is a “DAA” file, a binary file containing all aligned sequences and reported alignments.
- Then run tool daa-meganizer (or MEGAN)
  - to “meganize” the DAA file; performing taxonomic and functional analysis of all aligned sequences, and
  - the result of meganization is appended to the DAA file; no new file is created.
- A meganized DAA file can be opened in MEGAN.

# Taxonomy meganization

- Taxonomic binning uses
  - the NCBI taxonomy (Benson et al, 2005),
  - the GTDB taxonomy (Parks et al, 2018),
  - naive LCA algorithm for short reads (Huson et al, 2007),
  - interval-union LCA for long reads (Huson et al, 2018).



# Function meganization

- Functional binning uses e.g.
  - EggNOG (Powell et al, NAR 2014)
  - InterPro (Mitchell et al, NAR 2015)
  - SEED (Overbeek et al, NAR 2014)
  - KEGG (MEGAN UE only, Kanehisa & Goto, NAR 2000)
- Assignment uses the top-hit strategy (Huson, 2011)

# Meganization database

- A DAA file contains reference sequences and their accessions
- Meganization requires a mapping of accessions to taxonomic and functional classes
- Provided as a “MEGAN mapping database”  
`megan-map-tutorial.db`
- Here is the SQLITE schema:

```
CREATE TABLE mappings (Accession PRIMARY KEY, Taxonomy INT, GTDB  
INT, EGGNOG INT, INTERPRO2GO INT, SEED INT, EC INT);
```

- A typical entry:

EKP93748|867903||253||22932|501010007

# Step 2: Meganization

- Meganize one DIAMOND file:

```
~/megan/tools/daa-meganizer \
-i out/Alice00-1mio.daa \
-mdb megan-map-tutorial.db
```

- Meganize *all* DIAMOND files:

```
~megan/tools/daa-meganizer \
-i out/*.daa -mdb megan-map-tutorial.db
```

# Step 2: Meganization

- If you failed to meganize the 12 files, you can download the meganized files here:

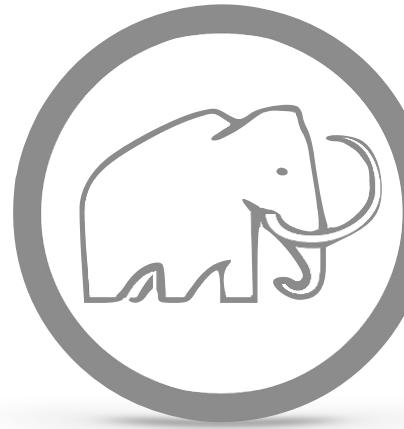
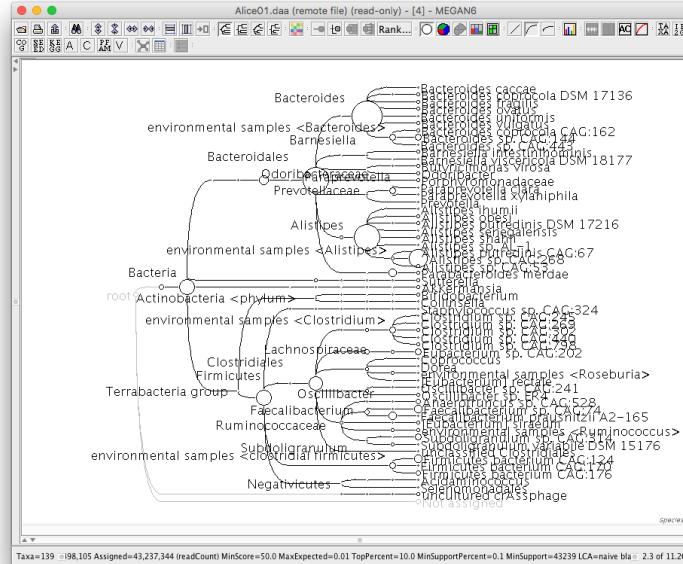
[https://software-ab.cs.uni-tuebingen.de/download/  
megan6/tutorial/meganizer-out.zip](https://software-ab.cs.uni-tuebingen.de/download/megan6/tutorial/meganizer-out.zip)

# Outline

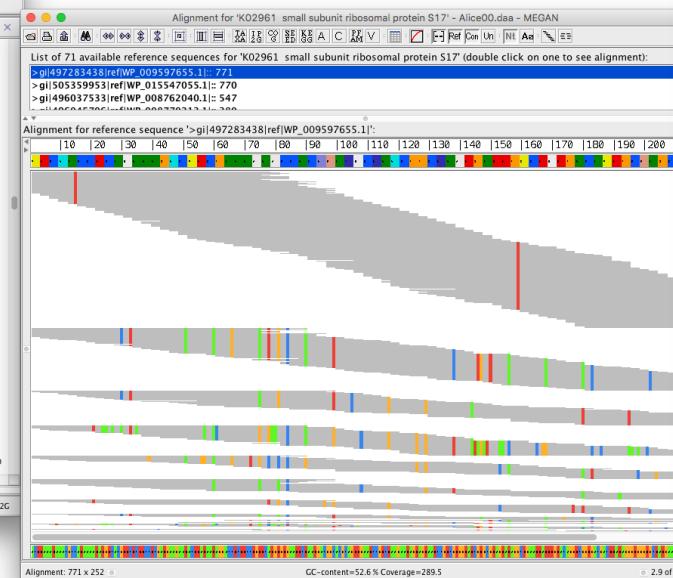
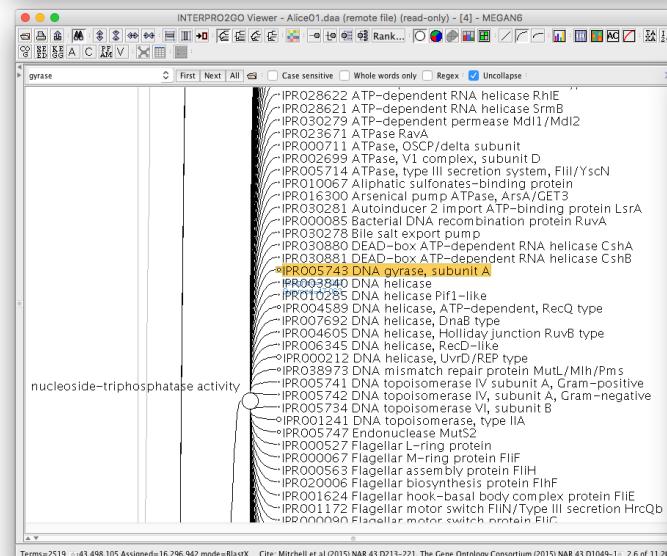
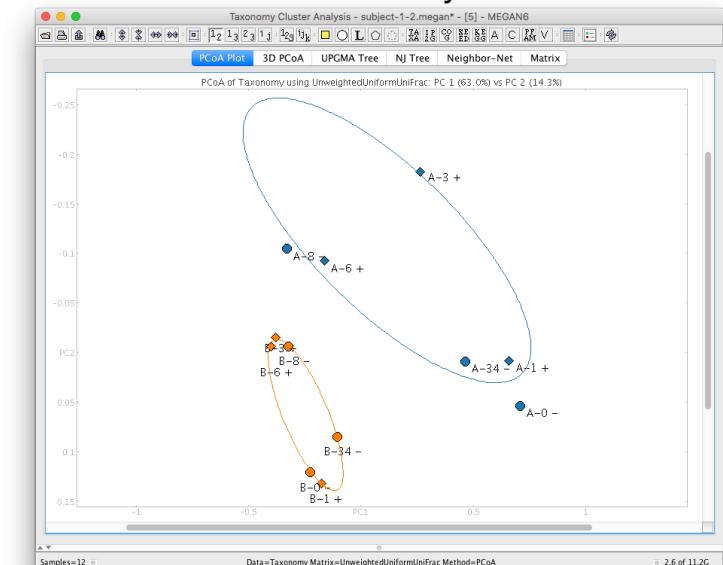
- Introduction to microbiome analysis
- Step 0: Software setup
- Step 1: DIAMOND alignment against protein database
- Step 2: MEGANization of reads and alignments
- Step 3: MEGAN interactive analysis

# Interactive MEGAN analysis

## Taxonomic content

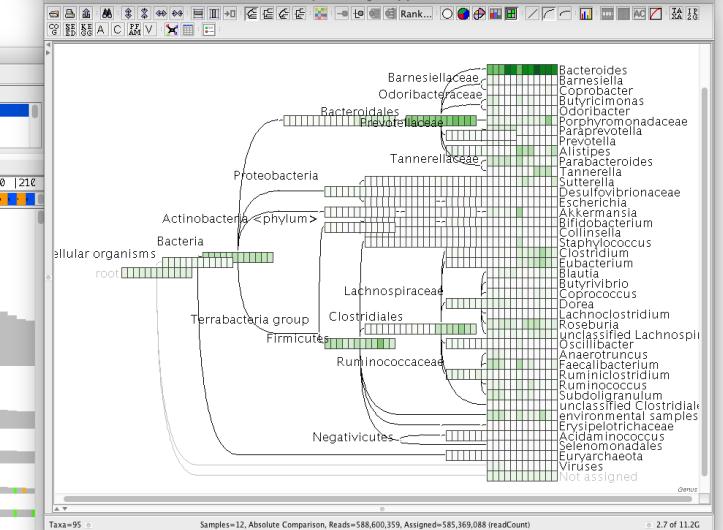


## PCoA analysis



## Functional content

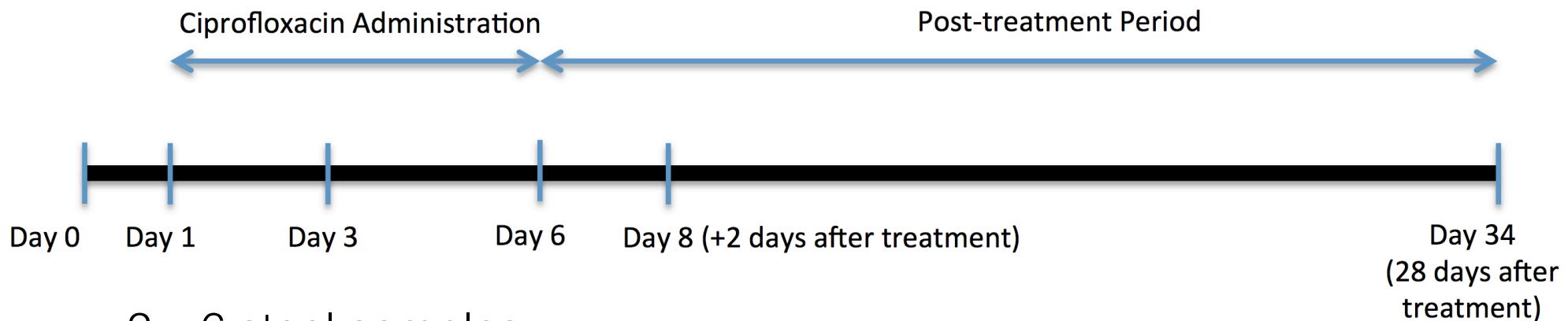
## Gene-centric alignment and assembly



## Comparative analysis

# ASARI- Antibiotic resistance pilot study

- Two volunteers, subject 1 and subject 2



- 2 x 6 stool samples
- Shotgun sequencing
  - ~60 million reads per sample (101 bp per read)
  - ~800 million reads in total
- Initial analysis: compare against NCBI-nr protein database

# Performance of DIAMOND+MEGAN

- 12 human gut samples, total 816 million HiSeq reads

Sample	Reads	DIAMOND (s)	Alignments	Aligned reads	Meganizer (s)
Alice 0	66 393 401	19 062	627 405 772	44 900 227	9 299
Alice 1	64 923 975	15 771	595 715 349	43 498 105	11 338
Alice 3	55 092 349	13 435	515 249 349	37 675 494	8 621
Alice 6	66 289 376	16 801	910 892 059	52 627 776	11 771
Alice 8	57 957 661	14 134	790 946 244	45 358 448	13 911
Alice 34	64 380 386	15 615	608 114 143	44 741 897	11 962
Bob 0	61 232 588	14 573	825 213 917	48 882 884	12 058
Bob 1	65 763 766	16 203	841 038 616	51 408 892	12 270
Bob 3	89 034 641	34 598	1 233 571 041	72 017 720	15 789
Bob 6	89 339 172	27 333	1 138 796 522	70 344 161	15 507
Bob 8	78 001 118	19 734	1 049 831 855	63 336 241	13 423
Bob 34	57 627 119	15 406	780 844 319	45568158	11 433
Total	816 035 552	222 665	9 917 619 186	620 360 003	Max: 15 789
Time		≈ 62 h			≈ 5 h

doi:10.1371/journal.pcbi.1004957.t001

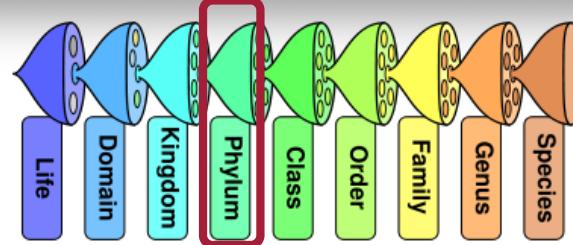
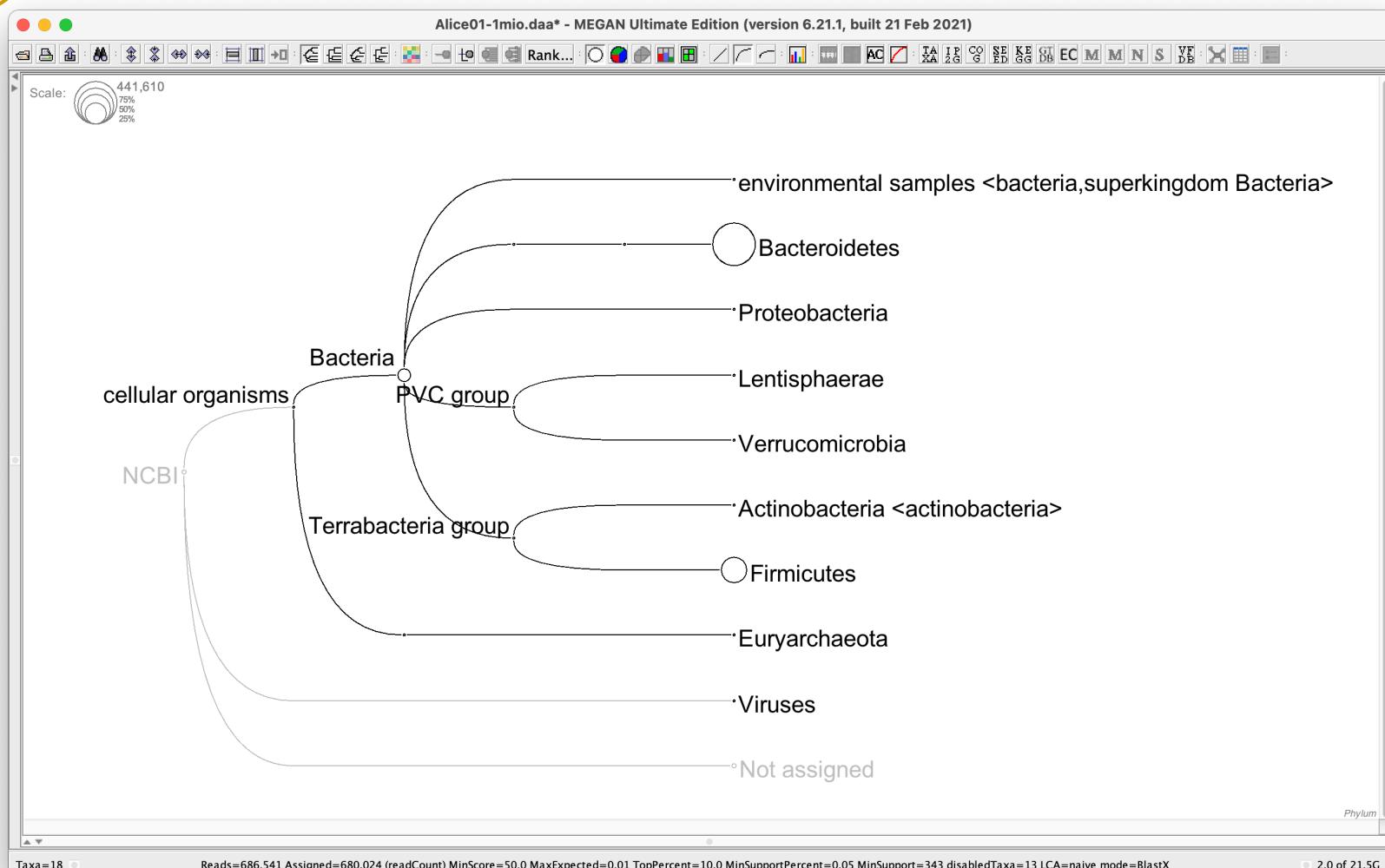
- Complete analysis in 62+5 hours on a single server



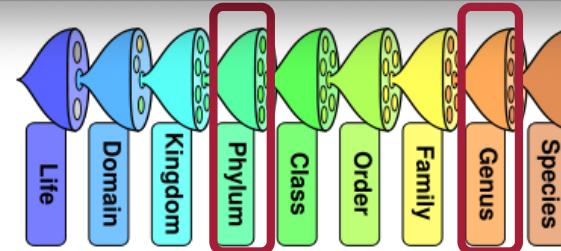
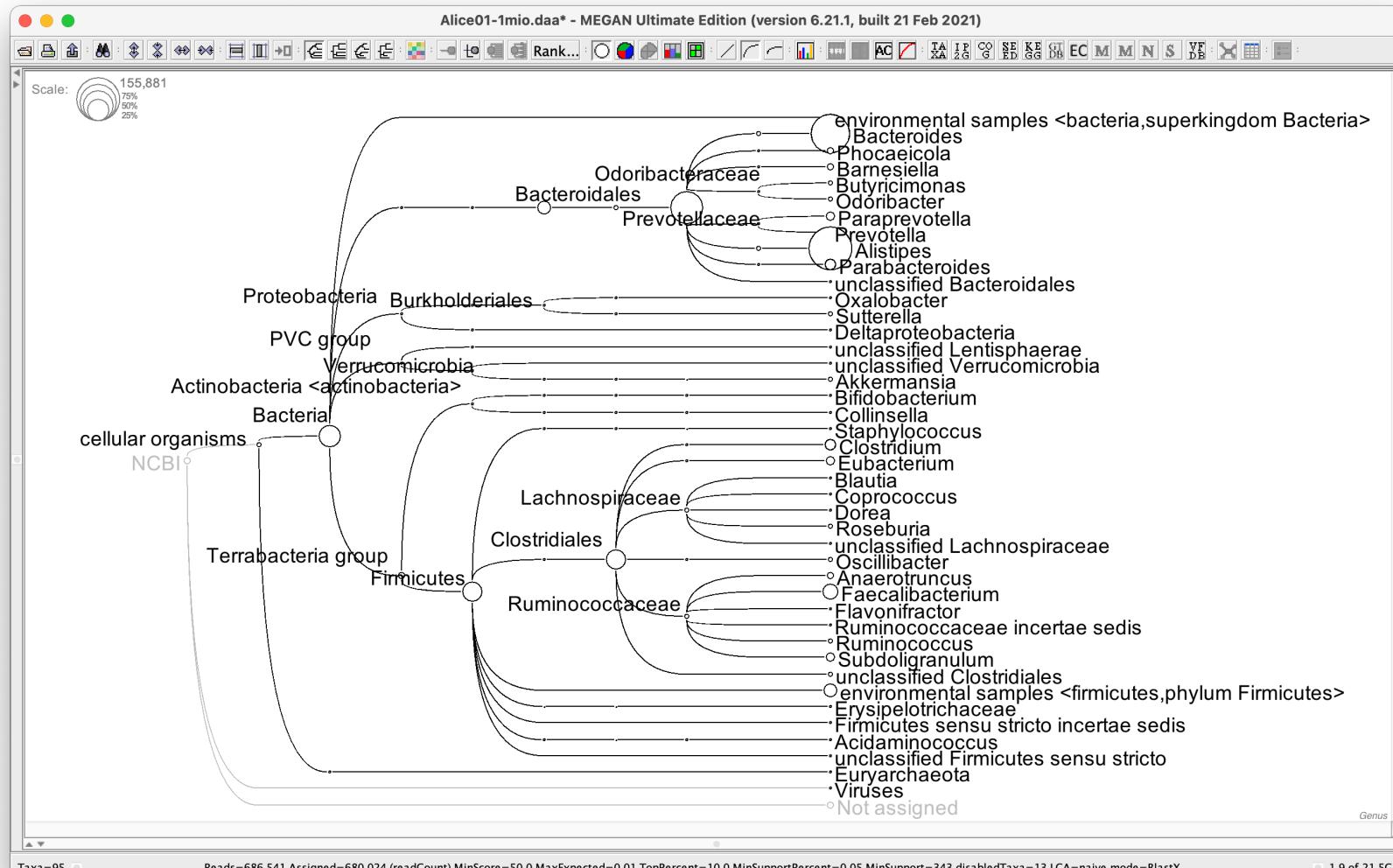


# Taxonomic content

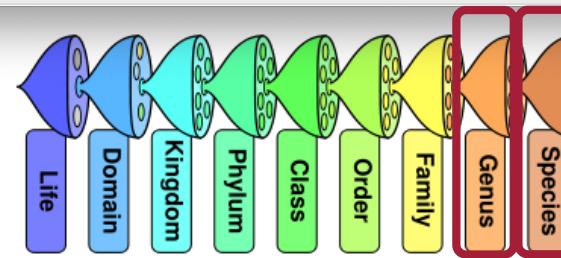
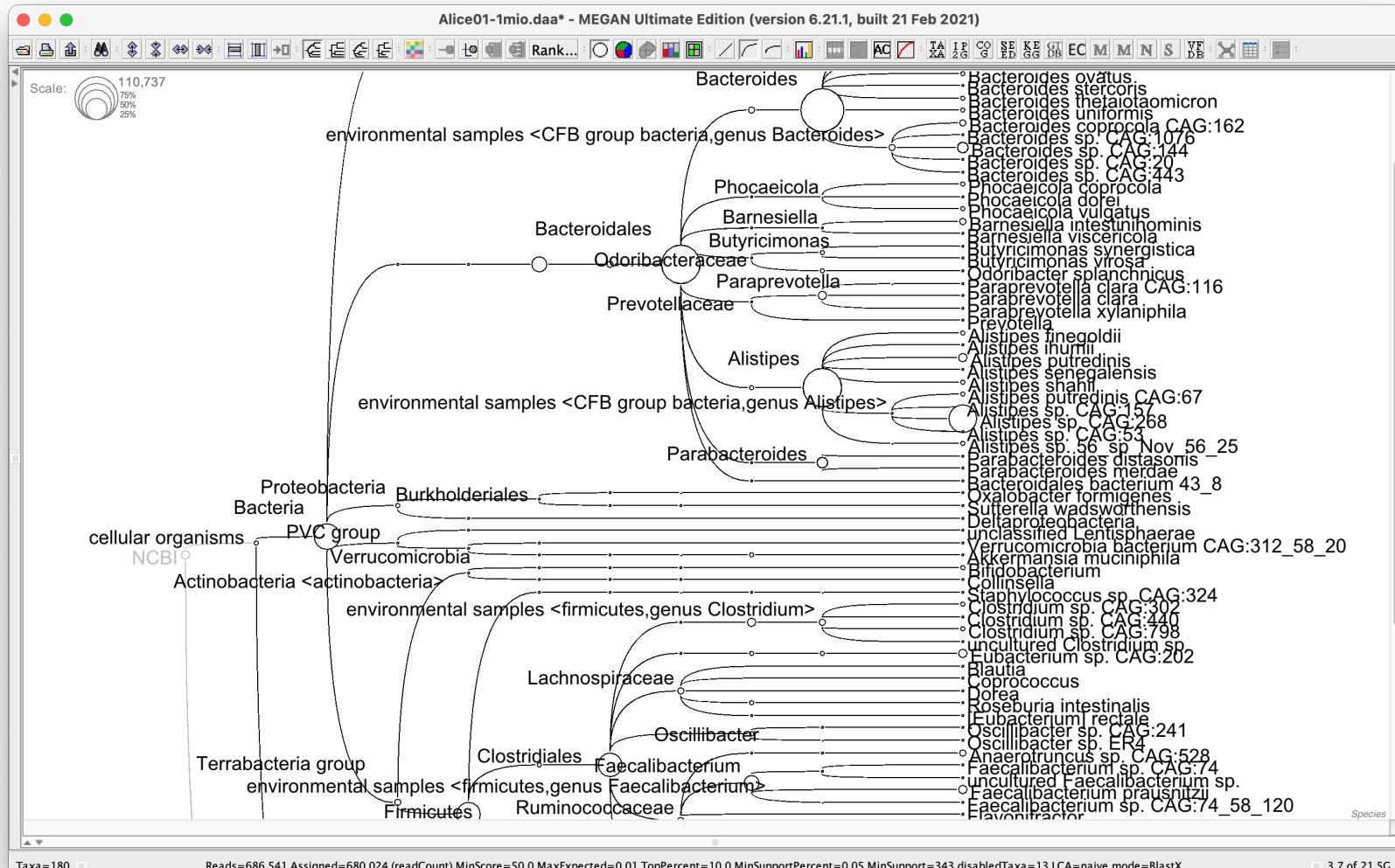
## ASARI human gut microbiome



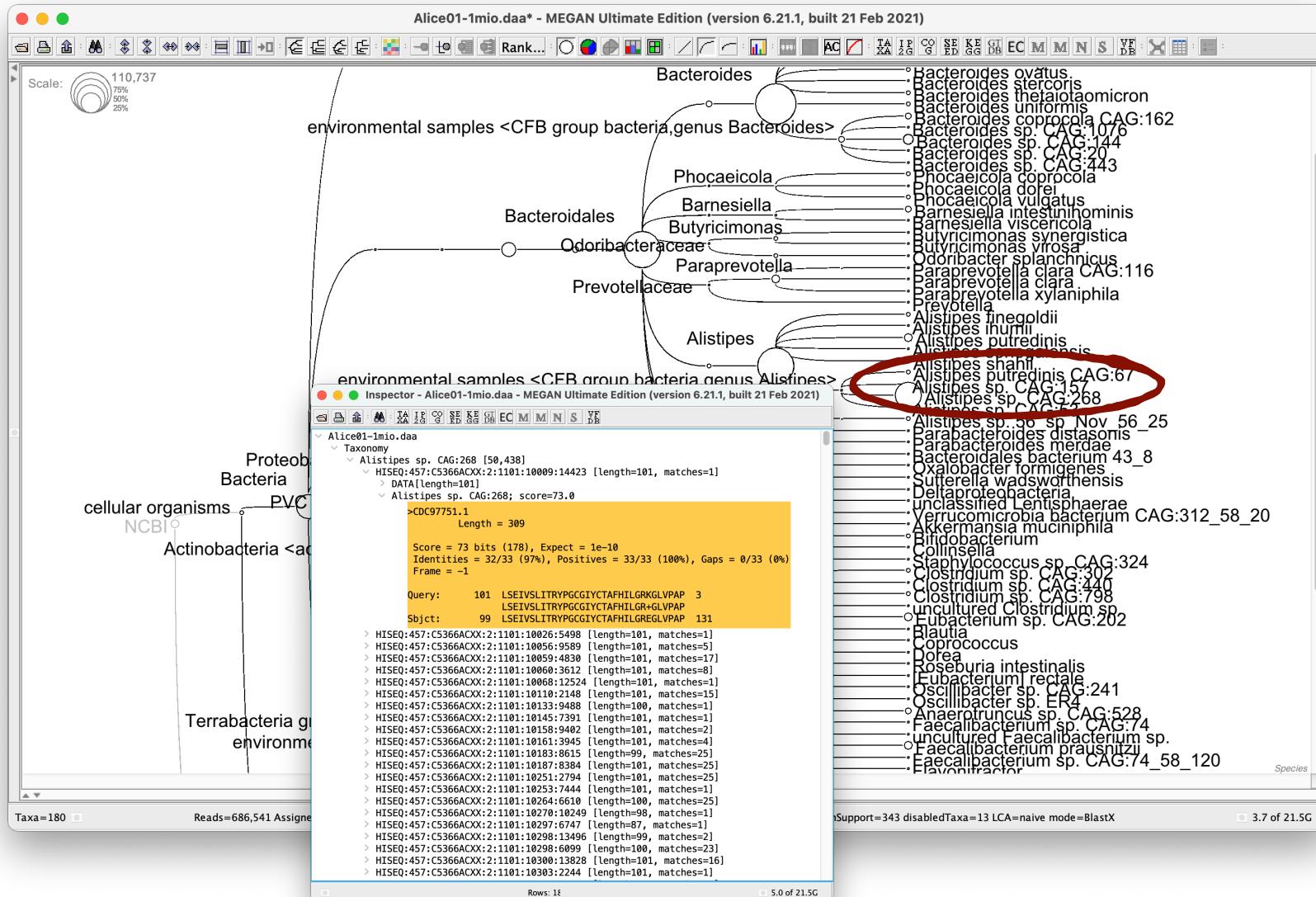
# Taxonomic content



# Taxonomic content



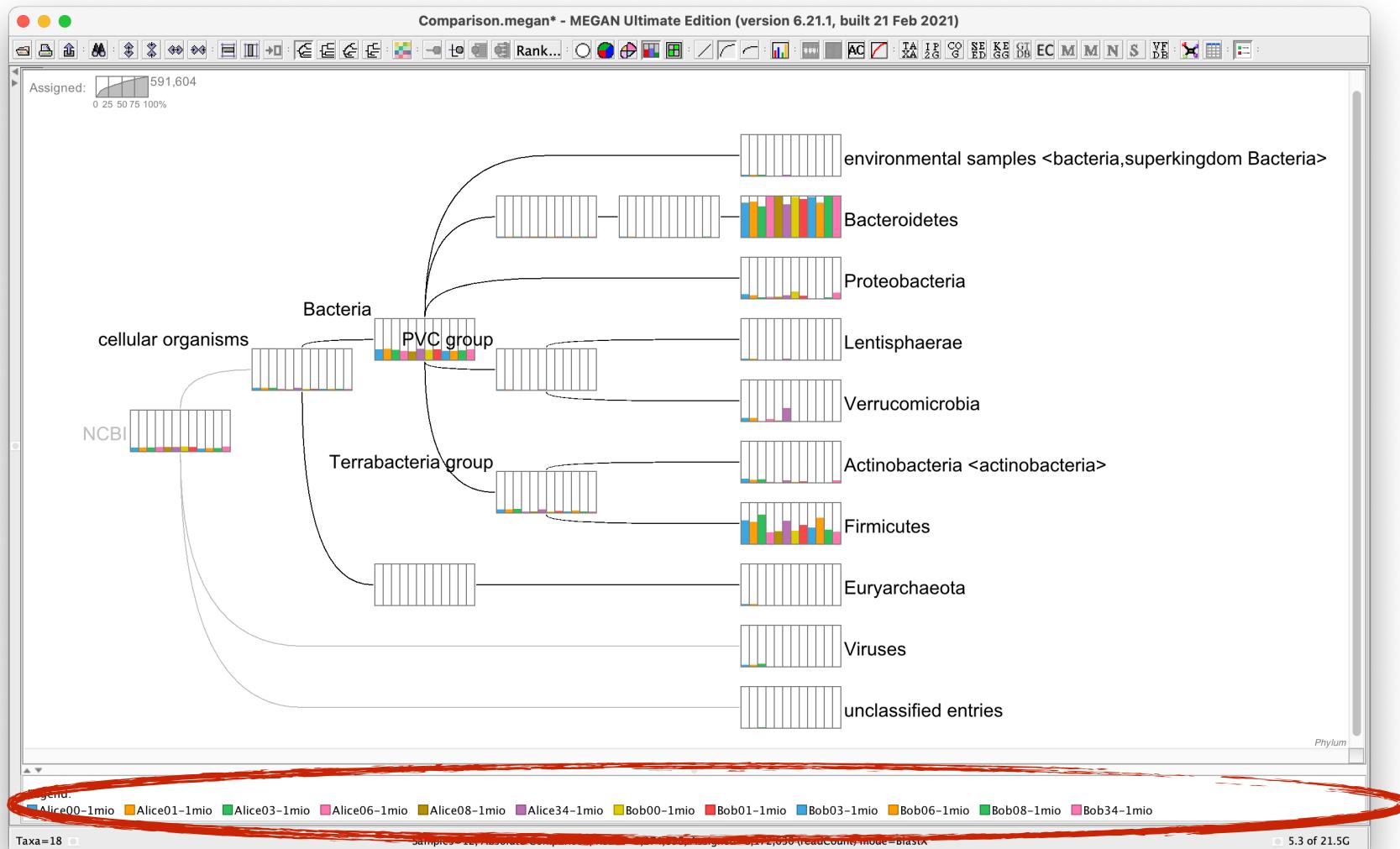
# Drill down to details...





Q3: How do they compare?

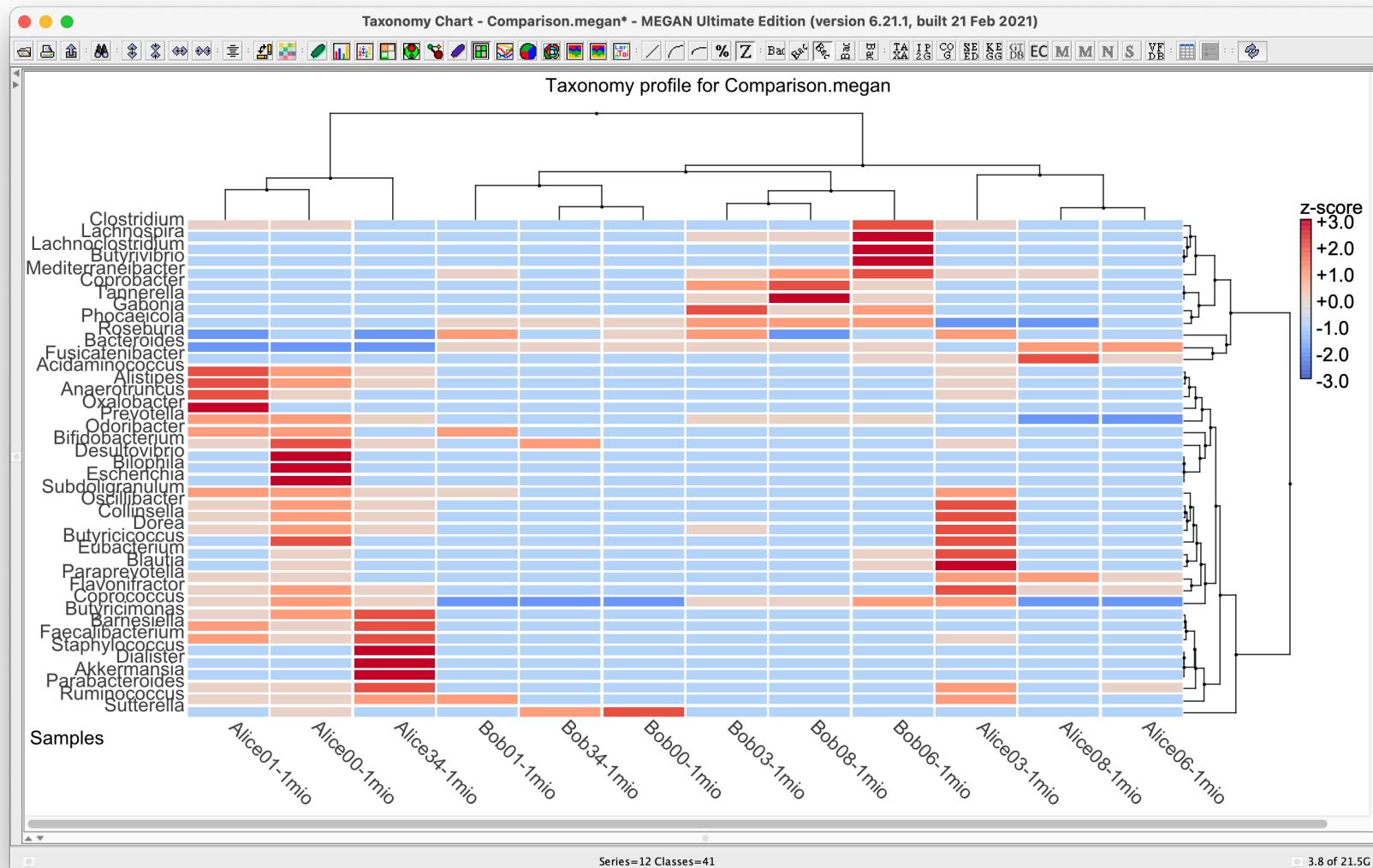
# Comparison



All 12 ASARI human gut samples together

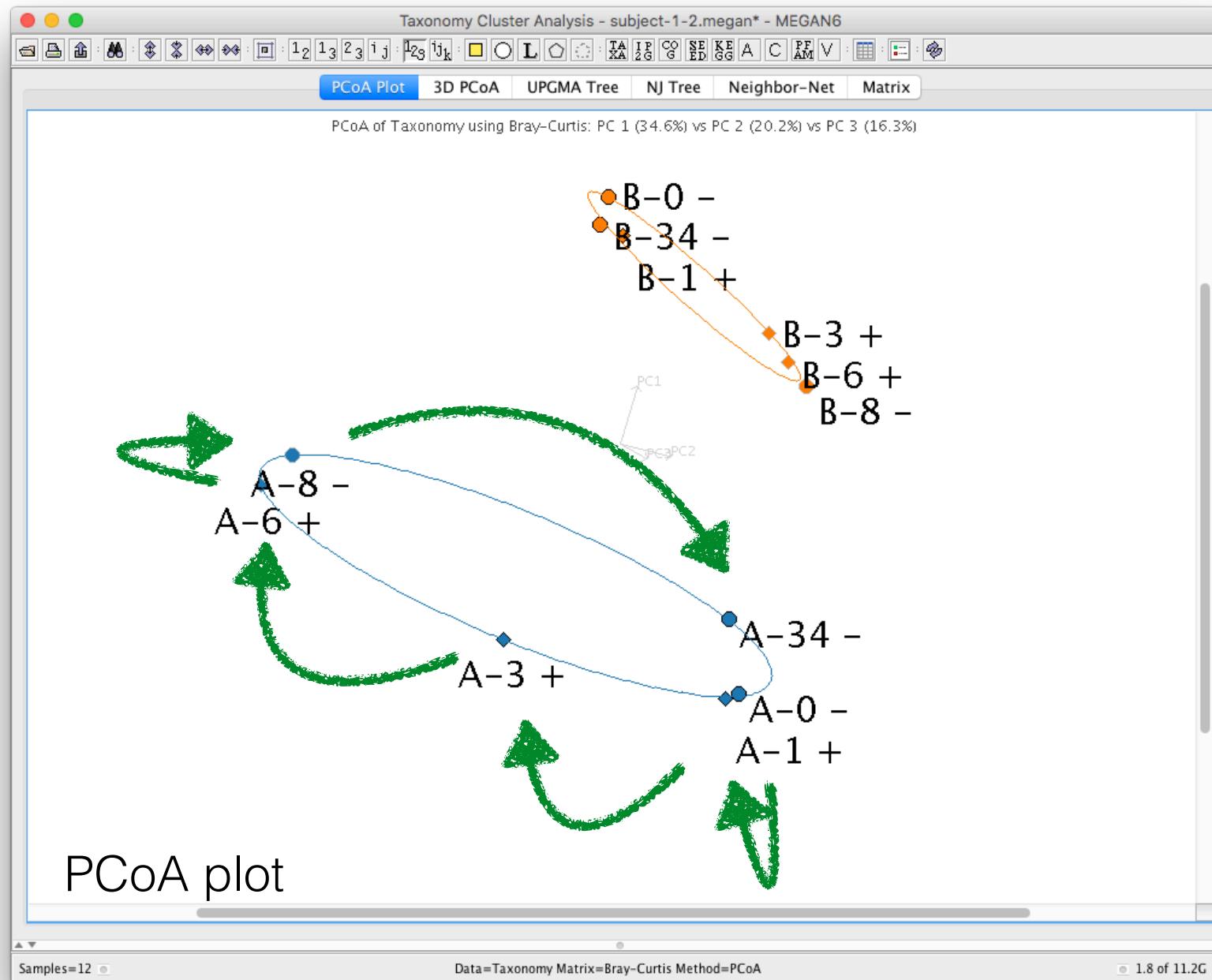


# Comparison

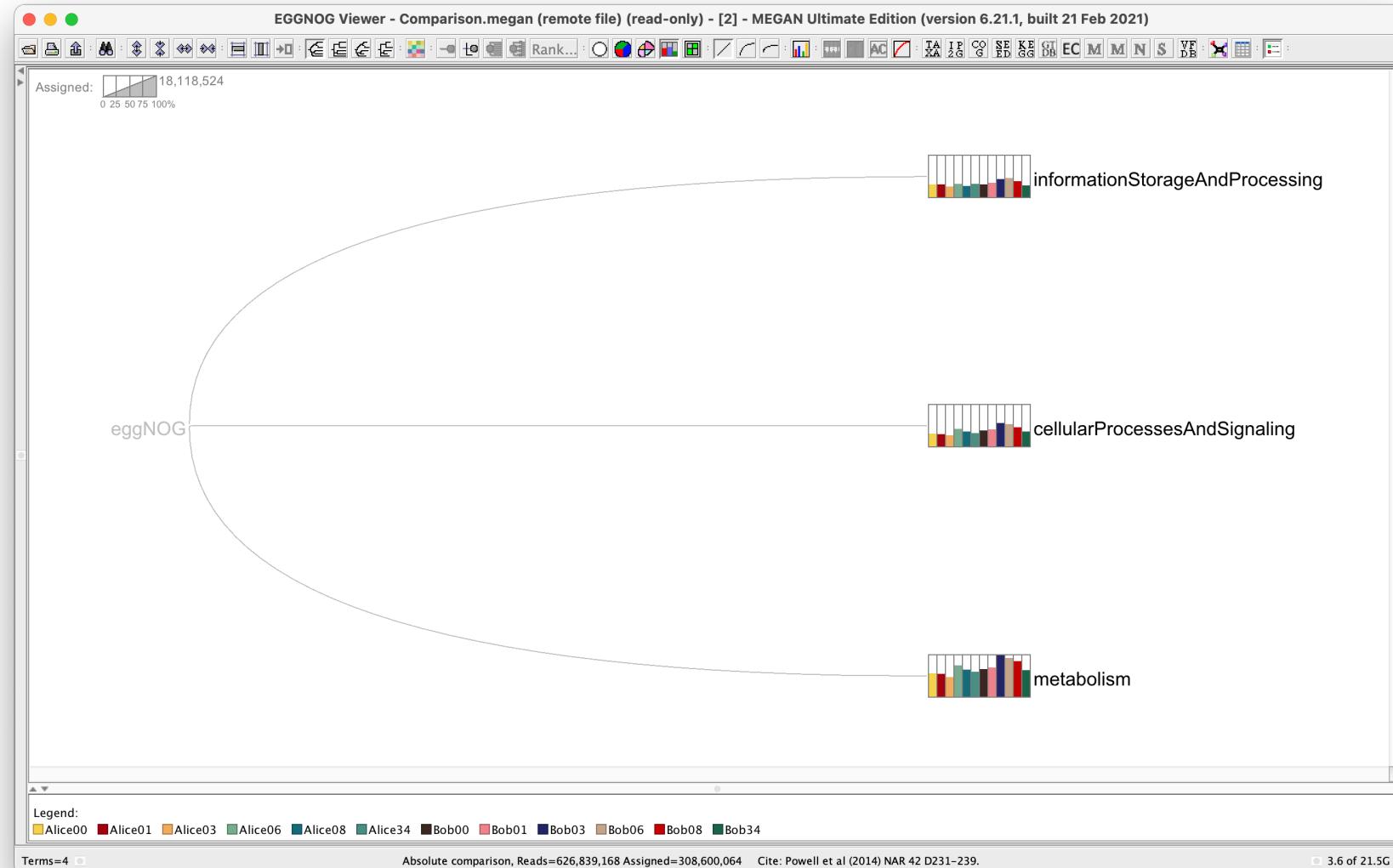


All 12 ASARI human gut samples together

# E.g.: Does the microbiome rebound?

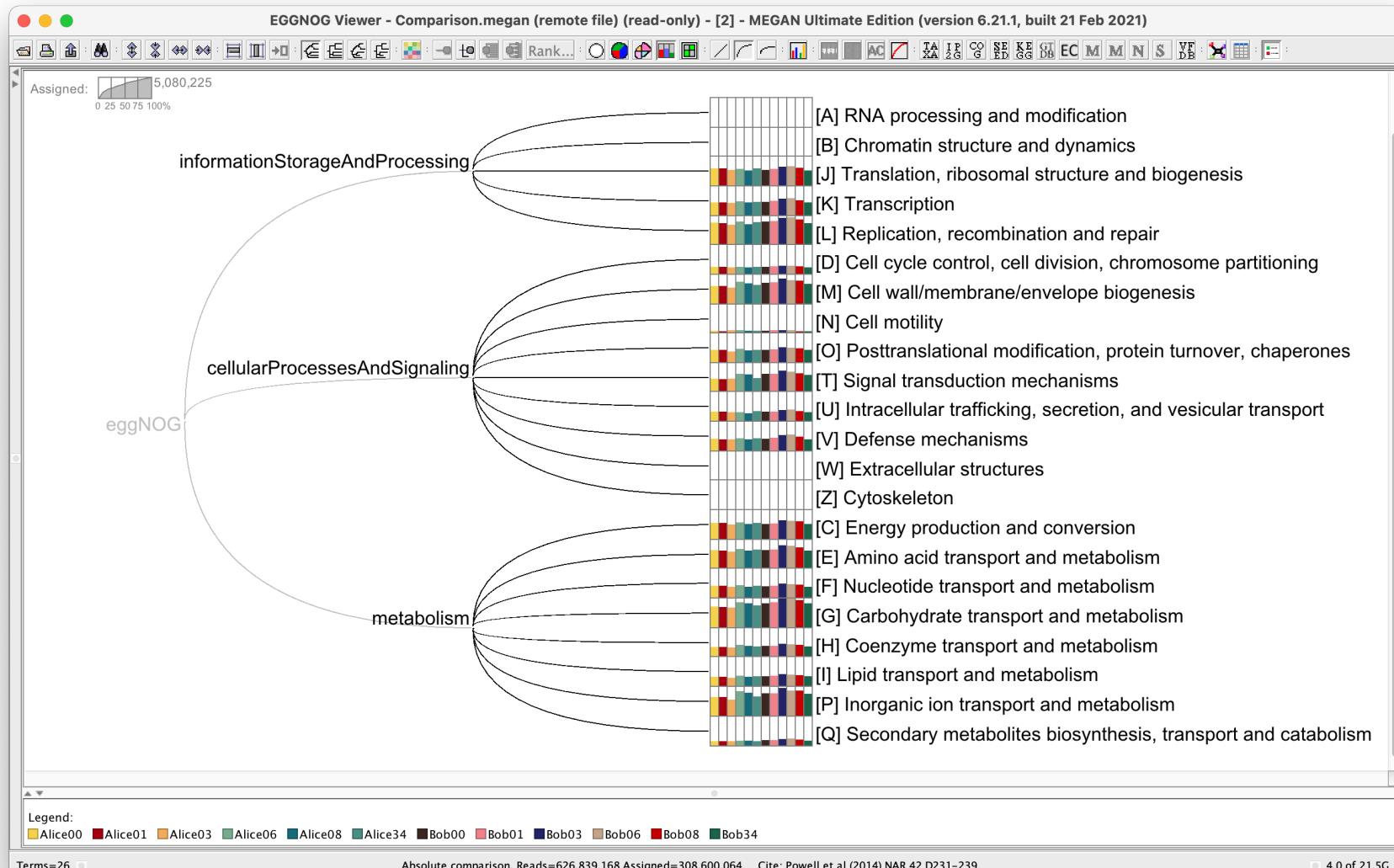


# Functional content



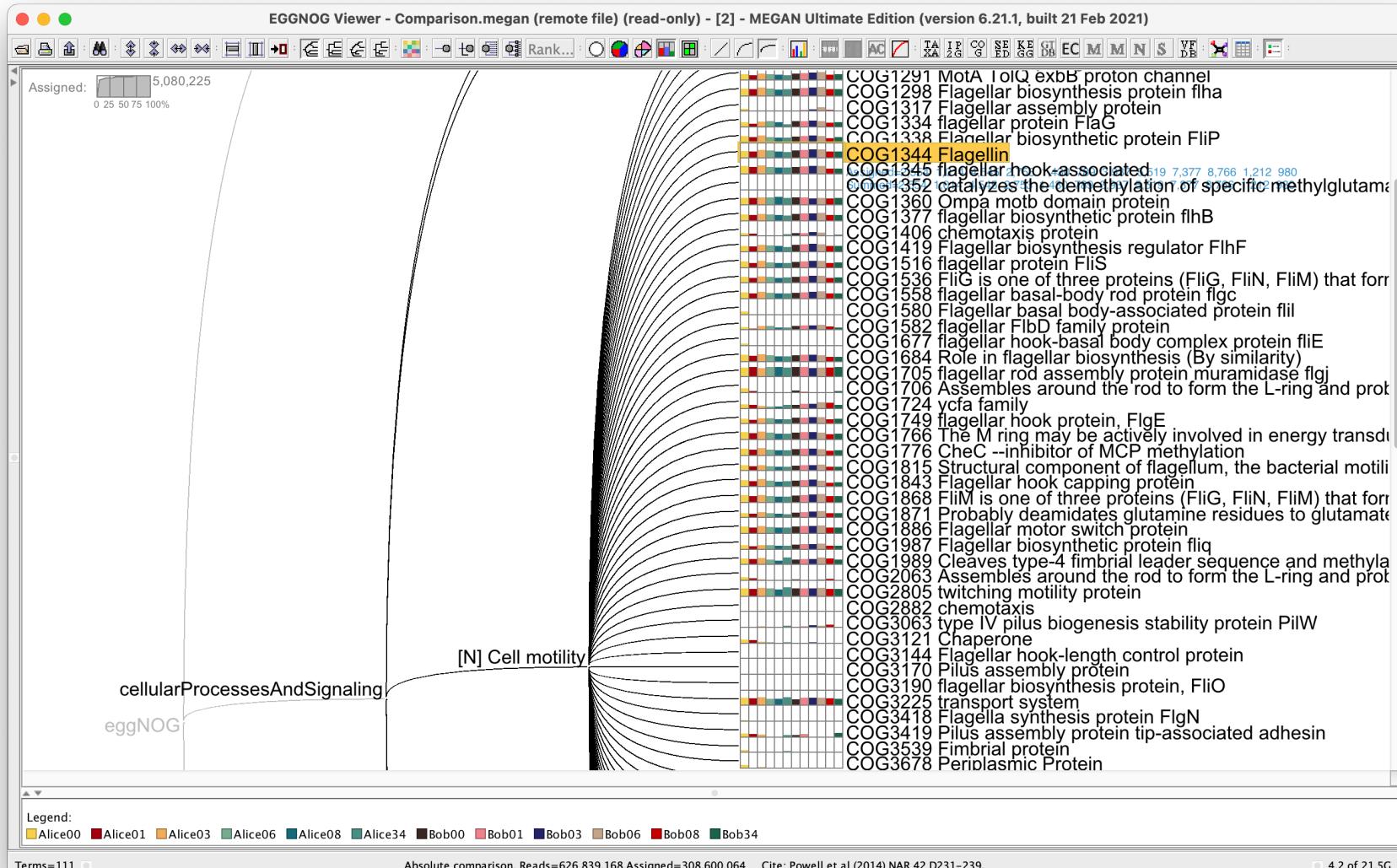
eggNOG classification  
(Powell et al, 2014)

# Functional content



eggNOG classification  
(Powell et al, 2014)

# Functional content



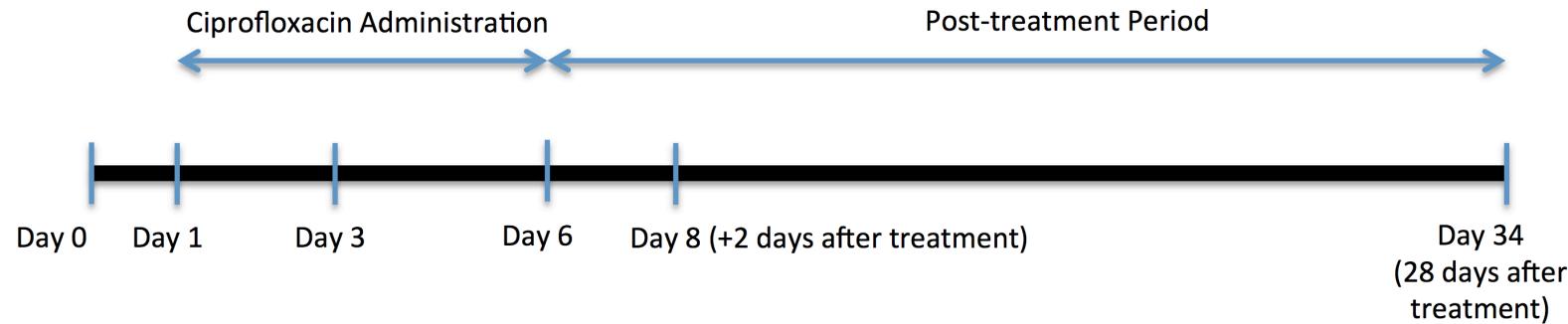
eggNOG classification  
(Powell et al, 2014)

# Step 3: MEGAN analysis

- Launch MEGAN by typing:  
`~ /megan/MEGAN`
- Open individual files with the File->Open... item
- Compare files using the File->Compare... item

# Alice and Bob- short reads

Alice and Bob, 6 time points each



- Each subsampled to 1 mio reads.
- `data/Alice00-1mio.fq.gz` etc

# Alice and Bob- short reads

- Tutorial tasks:
  - confirm that these are gut samples - should be dominated by Bacteroidota, Firmicutes and Proteobacteria.
  - Open all twelve files together in a comparison document and add the provided metadata.
  - Confirm that the taxonomic profiles of either subject changes during the course of antibiotics and then returns to a similar state after treatment.

# Enrichment reactor - long reads

- Nanopore reads from enrichment reactor:

Short report | [Open Access](#) | Published: 16 April 2019

## Annotated bacterial chromosomes from frame-shift-corrected long-read metagenomic data

Krithika Arumugam, [Caner Bağci](#), [Irina Bessarab](#), [Sina Beier](#), [Benjamin Buchfink](#), [Anna Górska](#), [Guanglei Qiu](#), [Daniel H. Huson](#) & [Rohan B. H. Williams](#) [✉](#)

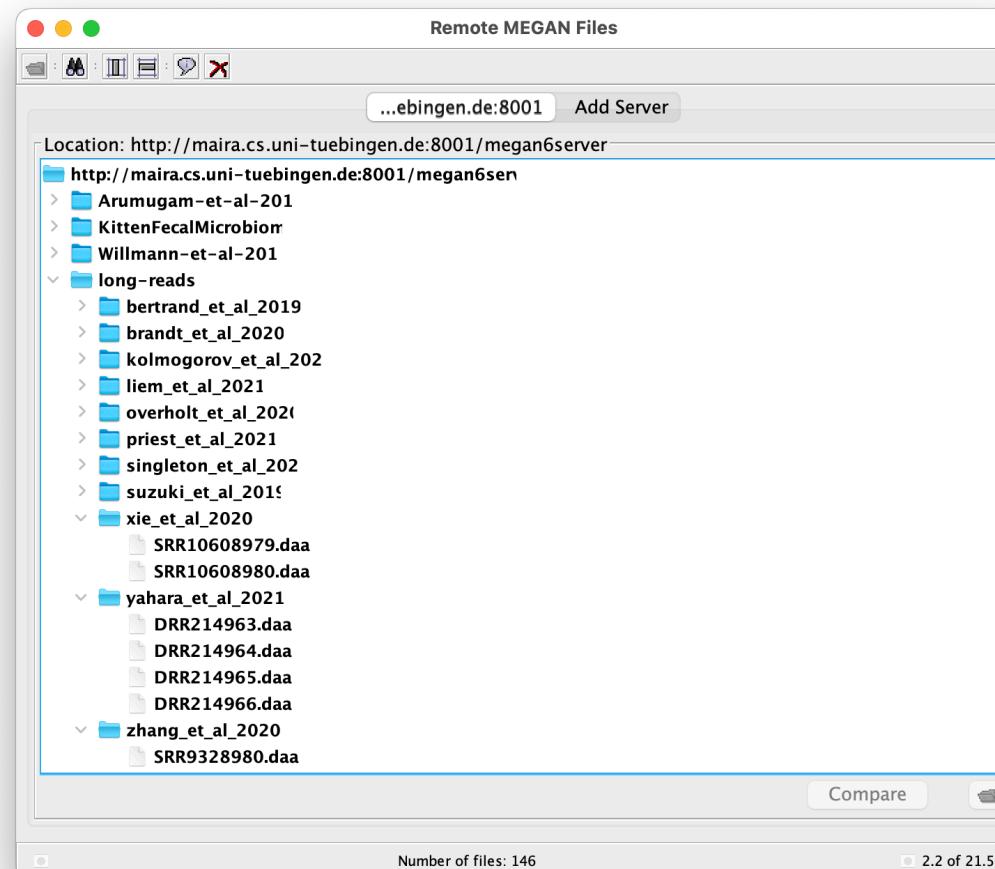
*Microbiome* 7, Article number: 61 (2019) | [Cite this article](#)



Krithika Arumugam

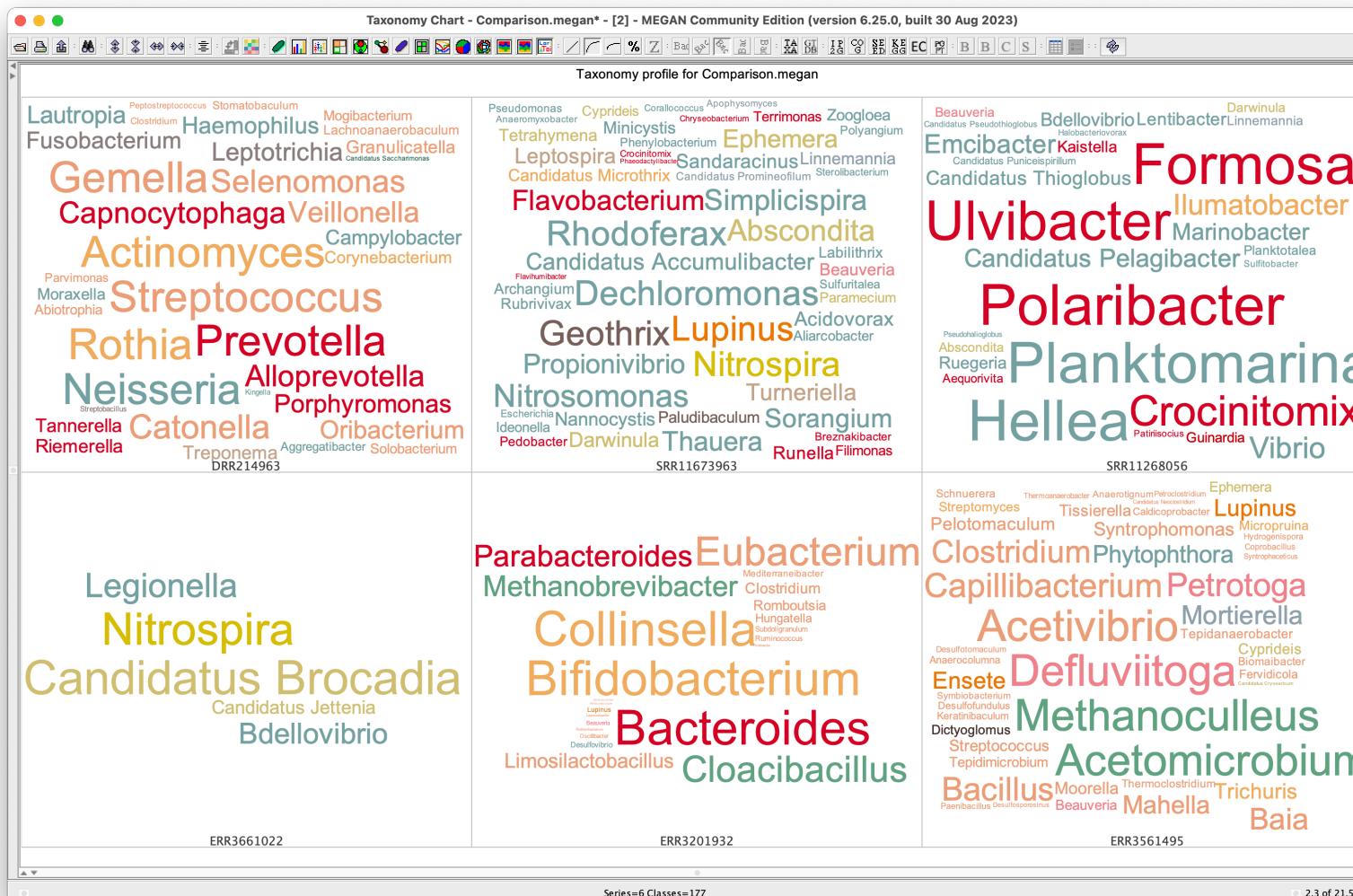
- Reads ~695,000, length ~9kb, total ~6Gb
- Unicycler assembly:  
long-reads/assembly.fa.gz

- MeganServer serves MEGAN files to the web
- The default server provides access to a number of published metagenome datasets:



# MeganServer

- Here are six long-read datasets:



- Can you match them to: biogas plant, ground water, human gut, oral, sea water, waste water ?

# Thank You!

## Joint work with:

- Caner, Baci, Banu Cetinkaya, Anupam Gautam, Timo Lucas, Sascha Patz, Patrick Wörz and Wenhuan Zeng

Tübingen

- Krithika Arumugam, Irina Bessarab & Rohan Williams

SCELSE/NUS Singapore

## Funding:

- Deutsche Forschungsgemeinschaft (MAIRA & BinAC)
- Life Sciences Institute at NUS
- NRF/MOE and NRF-EW, Singapore