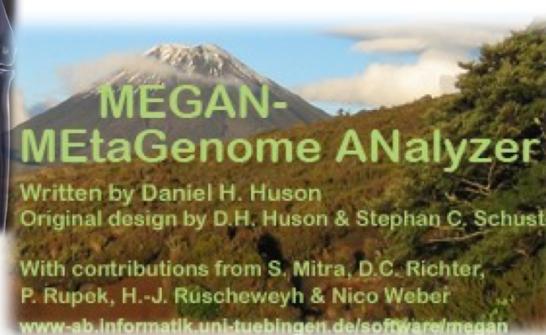


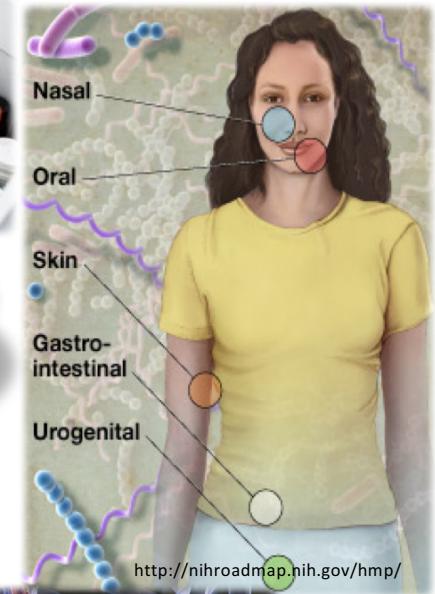
Metagenomics using MEGAN6

Who is out there?
What are they doing?
How do they compare?

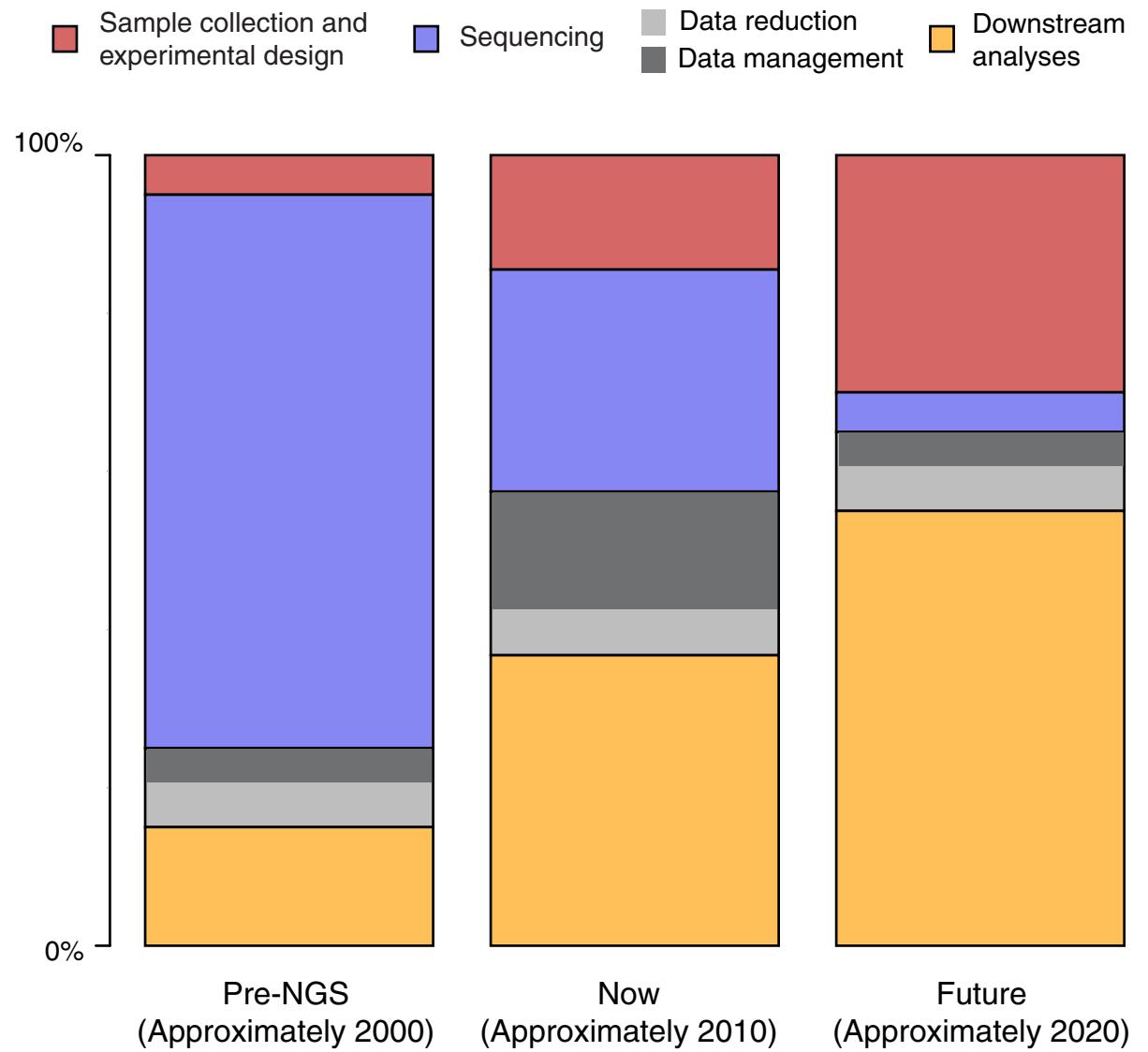
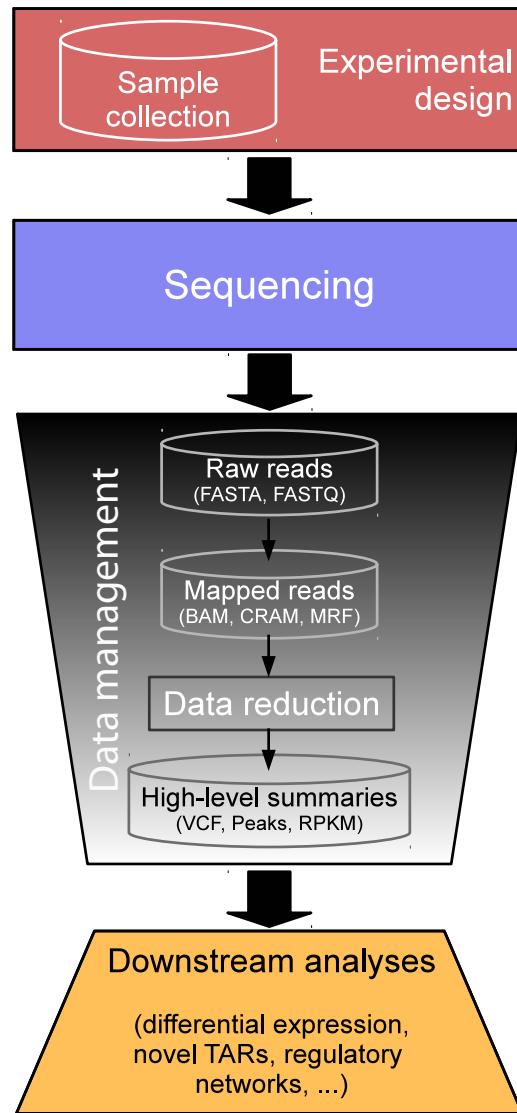




Metagenomics

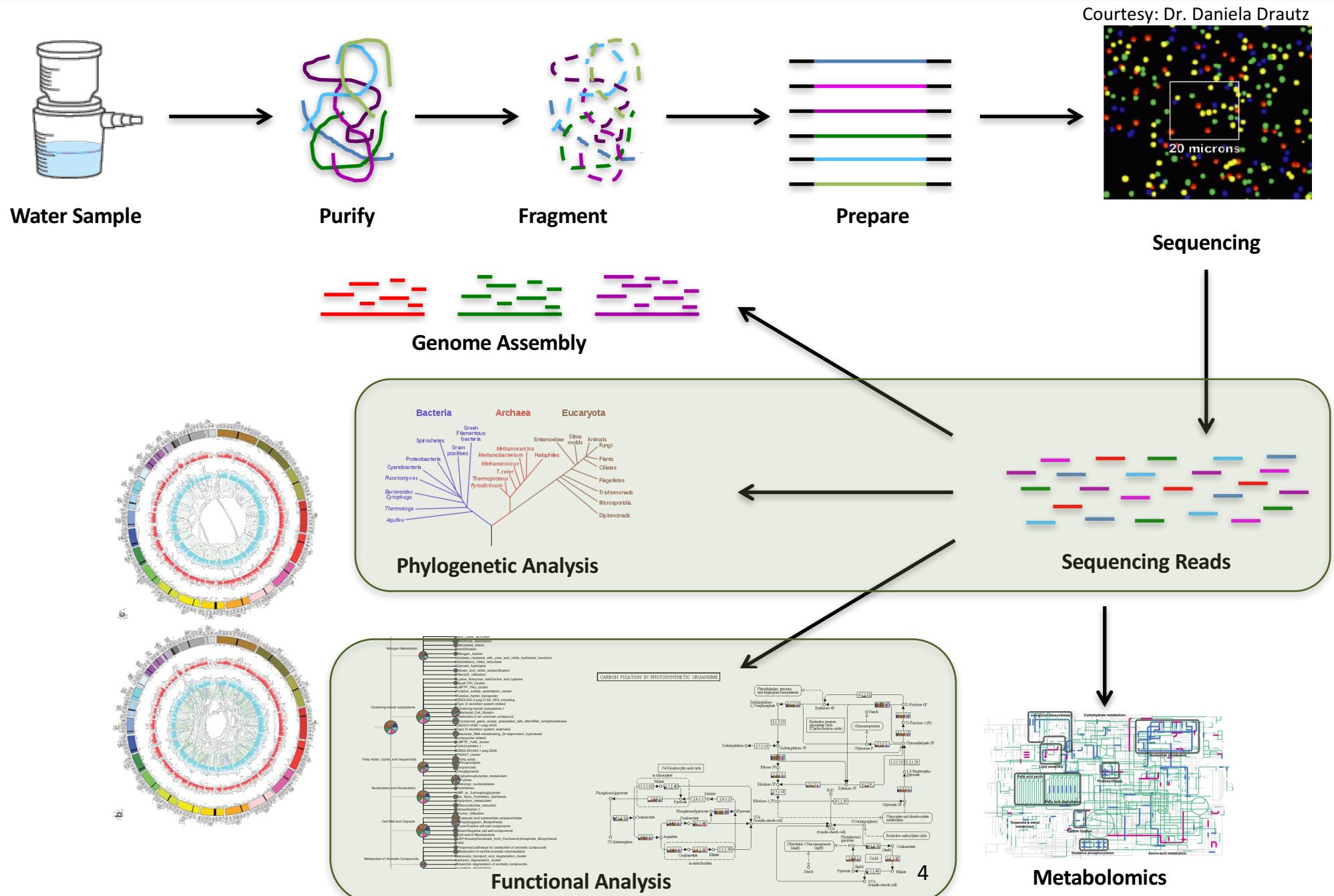


Sequencing ... Bottleneck?



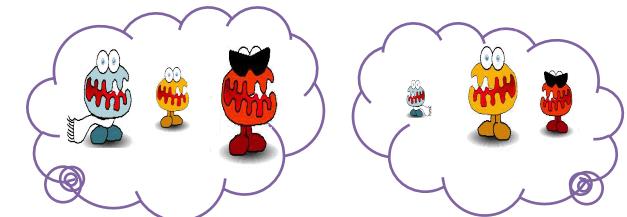
Sboner et al. Genome Biology 2011 12:125 doi:10.1186/gb-2011-12-8-125

Typical meta-omics study



Three Basic Computational Questions

- Who is out there?
 - Types of organisms
 - In what proportions?
- What are they doing?
 - Types of genes
 - Which metabolic pathways?
 - In what proportions?
- How do different samples compare?
 - Pairwise and multiple comparisons
 - Correlations with environmental parameters?
- Serve to answer biological or medical questions



Who is Out There?

Taxonomic analysis

Two main approaches:

- Targeted sequencing:
 - Sequence a specific gene, usually 16S rRNA, and place reads into a reference phylogeny
- Metagenome sequencing:
 - Randomly sequence DNA (or RNA) and then place reads into the NCBI taxonomy based on similarity to reference sequences

What They are Doing?

Functional Analysis

- A **subsystem** is a set of functional roles that implement a specific biological process or structural complex
- Genes are assigned to functional roles in subsystems using SEED
- Analysis of metabolic pathways using KEGG orthologies

Overbeek et al., Nucleic Acids Res 33(17), 2005

Kanehisa et al., Nucleic Acids Res 28, 2000



MEGAN – MEtaGenome ANalyzer



Huson et al, 2011

The screenshot displays three main windows of the MEGAN software:

- Phylogenetic Tree Window:** A tree diagram showing the taxonomic distribution of reads. The root branches into "cellular organisms" and "dsDNA viruses, no". "cellular organisms" further branches into "Bacteria" and "Archaea". "Bacteria" includes nodes for "Proteobacteria", "Bacteroidetes/Chlorobi group", "Spirochaetes", "Firmicutes", "Fusobacteria", "Cyanobacteria", "Actinobacteria", "Tenericutes", "Chlamydiae/Verr", and "Eukaryota". "Archaea" includes "Not assigned" and "No hits".
- Sequence Analysis Window:** Titled "Inspector - new.rma - MEGAN". It shows a detailed view of a sequence alignment between "Fusobacteria" and "Deinococcus radiodurans R1". The sequence score is 90.9 bits (224), Expect = 4e-17, Identities = 57/149 (38%), Positives = 78/149 (52%). The frame is -1. The sequence details include:
 - Query: 515 IAIDPS*VSCKAGKTTAHIGRFWGSACASVAKHGLEILGIAVIVADIRDAMMLRAVQTINST 3 +AID S KAG+ TA+G FW+GC+ + G+E A+ID R A+ + QL + Sbjct: 97 LAIDASFHRRKAGQHTAHLSFWNGCARTERGIEQSCLALDVQHRQALTVDVRQTLTG 1
 - Query: 335 ELEGKKFTLNQWVLSQLVLTQRTDLKITSLLVADAASFVSLPVFVEGLKEIGFSLSRLRN 1 E + VL RT + +VAD ++ P VE + G ISRL N Sbjct: 157 EAPTRLQEAXQLDQVLLDLRTVQLDLAAVVADGNYAKEPIVETVTGHLGPISLRLRN 2
 - Query: 155 AVLYYYIYEGPRTGKGRGPFTKGKGDIFSN 69 A L +Y G +RGR K DGR+DFS+ Sbjct: 217 ANLNDLYTGEHIPARRGKKKFDGKVDFSD 245
- Functional Analysis Window:** Titled "Functional Analysis (SEED) - new.rma* - MEGAN". It shows a hierarchical tree of SEED terms under "Secondary Metabolism". Some terms listed include: Carbohydrates, Cofactors, Vitamins, Prost, Virulence, Sulfur Metabolism, Stress Response, Protein Metabolism, RNA Metabolism, Motility and Chemotaxis, Biosynthesis of phenylpro, Bacterial cytostatics, differ, Plant Hormones, Biologically active compou, Aromatic amino acids and, Amino Acids and Derivat, Nitrogen Metabolism, Clustering-based subsys, Fatty Acids, Lipids, and Is, Nucleosides and Nucleoti, Cell Wall and Capsule, Metabolism of Aromatic C, DNA Metabolism, Phages, Prophages, Trans, Membrane Transport, Respiration, Miscellaneous, Cell Division and Cell Cyc, Dormancy and Sporulation, Regulation and Cell Signa, Phosphorus Metabolism, Potassium metabolism.

Bottom status bar: Taxa=18 ~ Reads=7353 Assigned=7063 MinScore=35.0 TopPercent=10.0 MinSupport=5 disabledTaxa=9 ~ 1215 of 2796M

Bottom status bar: SEED terms=34 ~ Reads=7353 Assigned=3755 ~ 1137 of 2796M

- Interactive tool for metagenomic analysis



MEGAN6

The most powerful interactive microbiome analysis tool

[Overview](#) [What's New](#) [Features](#) [Getting Started](#) [Download](#) [Buy & Upgrade](#)



The most powerful interactive microbiome analysis tool
Analyse metagenome, metatranscriptome and amplicon sequences from multiple sources

[**Download MEGAN6**](#)

[**What's new in MEGAN6**](#)

Microbiome analysis using a single application

MEGAN6 is a comprehensive toolbox for interactively analyzing microbiome data. All the interactive tools you need in one application.

- Taxonomic analysis using the NCBI taxonomy or a customized taxonomy such as SILVA
- Functional analysis using InterPro2GO, SEED, eggNOG or KEGG
- Bar charts, word clouds, Voronoi tree maps and many other charts
- PCoA, clustering and networks
- Supports metadata
- MEGAN parses many different types of input

Why use MEGAN6?

The software is:

1. Easy to use. MEGAN6 is a single application and all features are available through menus, toolbars and graphics. No scripting skills required.
2. Powerful. MEGAN6 allows you to work with hundreds of samples containing hundreds of millions of sequencing reads. Blast-like analysis can be performed using DIAMOND.
3. Comprehensive. MEGAN6 offers a large range of analysis tools, and is under active development.

Choose your Edition

Community Edition

Open Source version of MEGAN6

Ultimate Edition

Premium version of MEGAN6
Additional features,
tools and support

[**Download Community**](#)

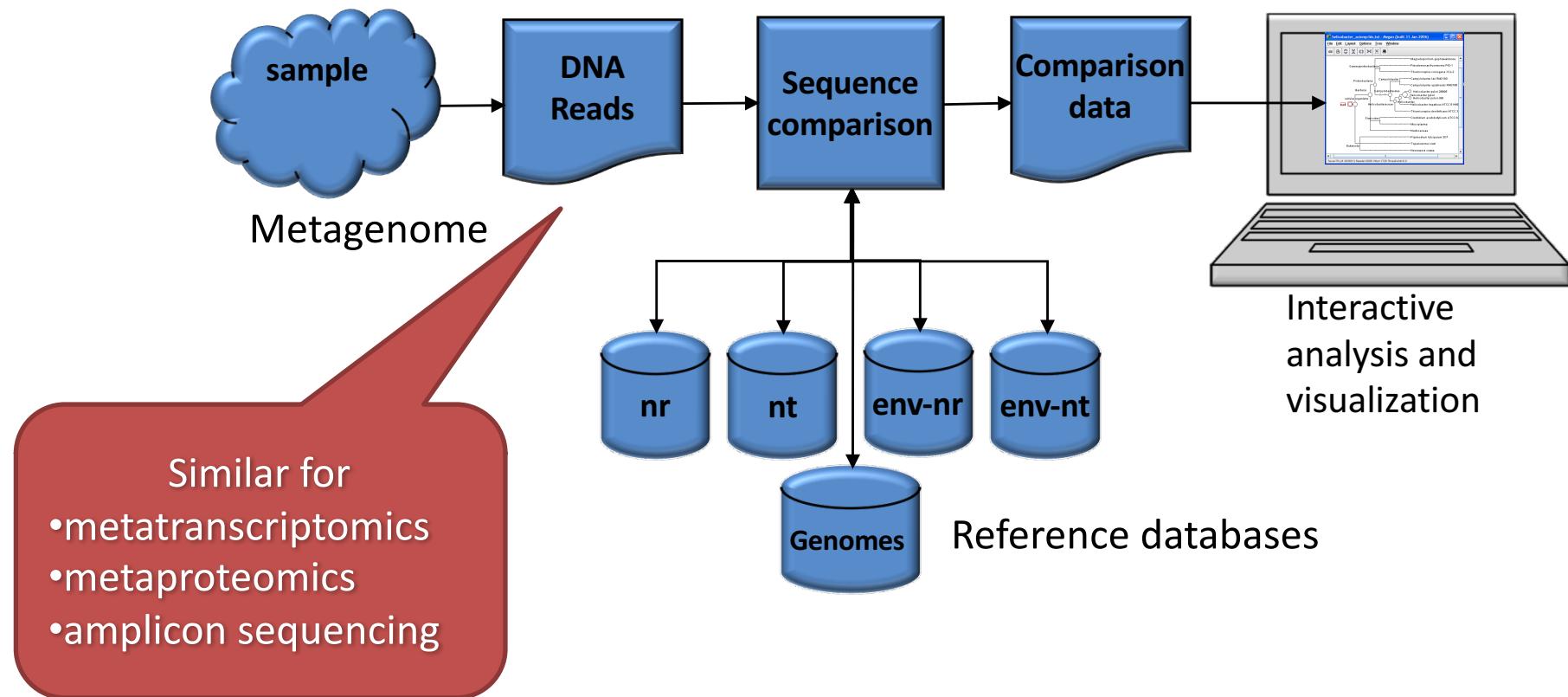
Free download

[**Download Ultimate**](#)

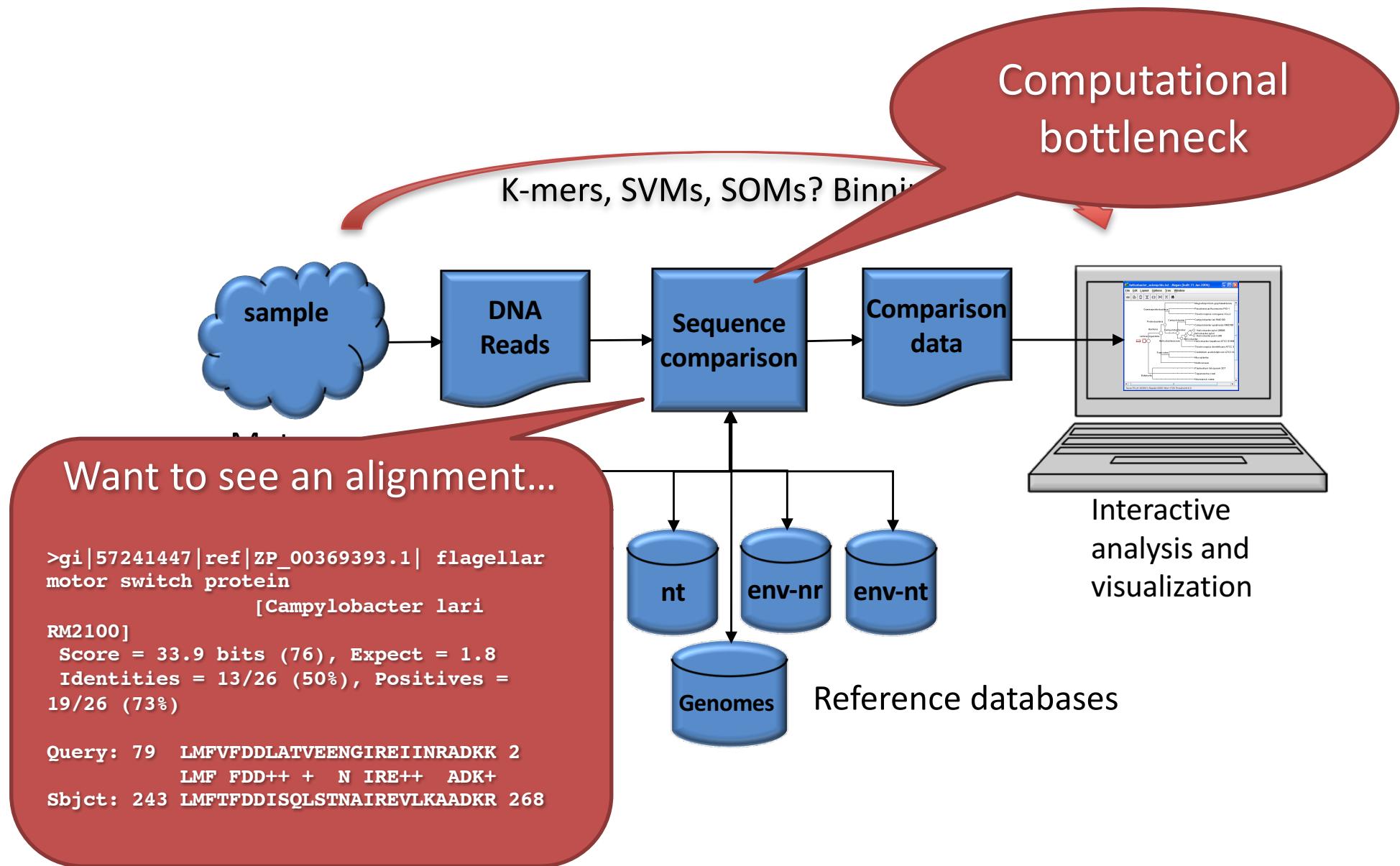
Free 30-day trial



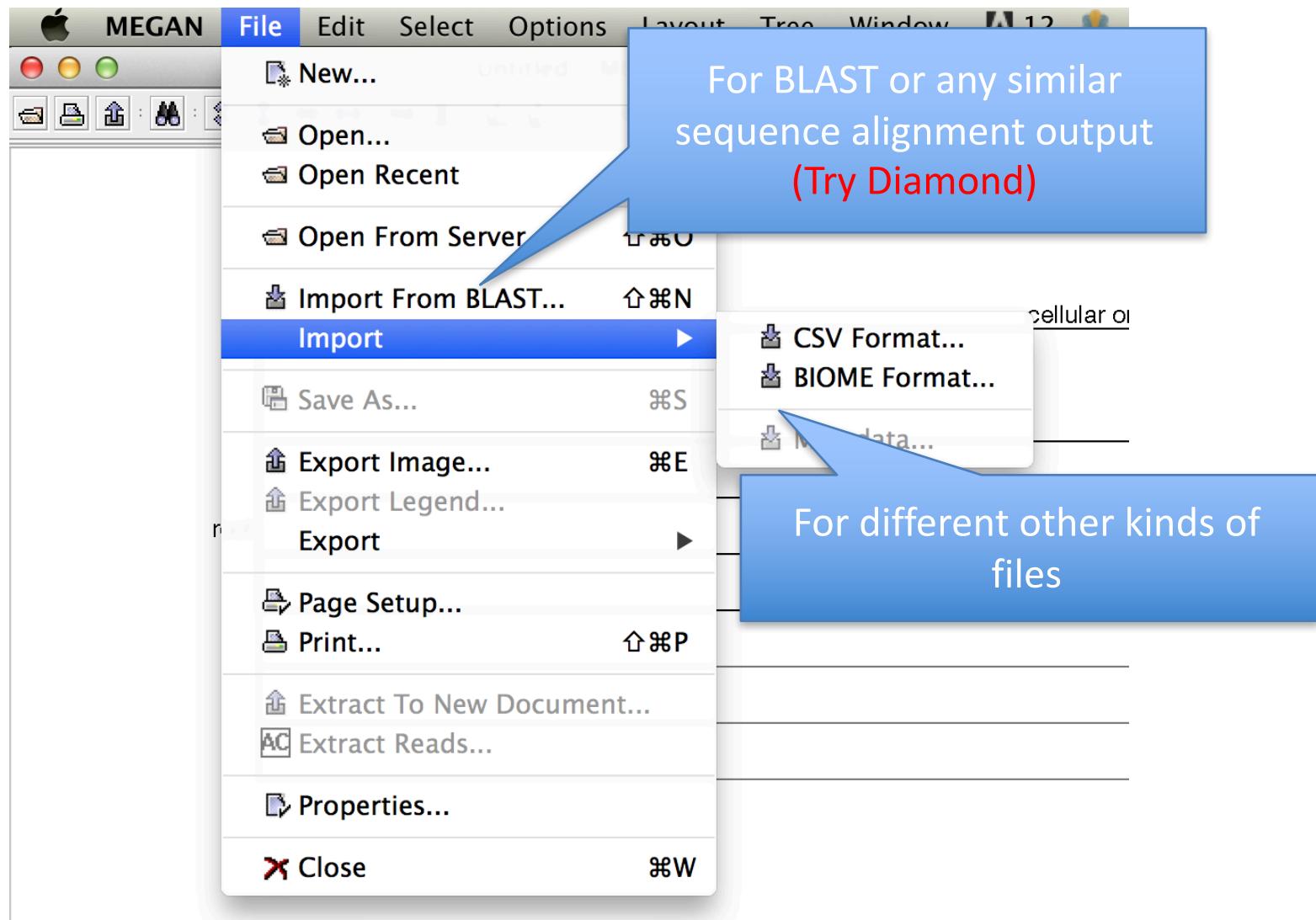
Metagenomics Pipeline



Metagenomics Pipeline

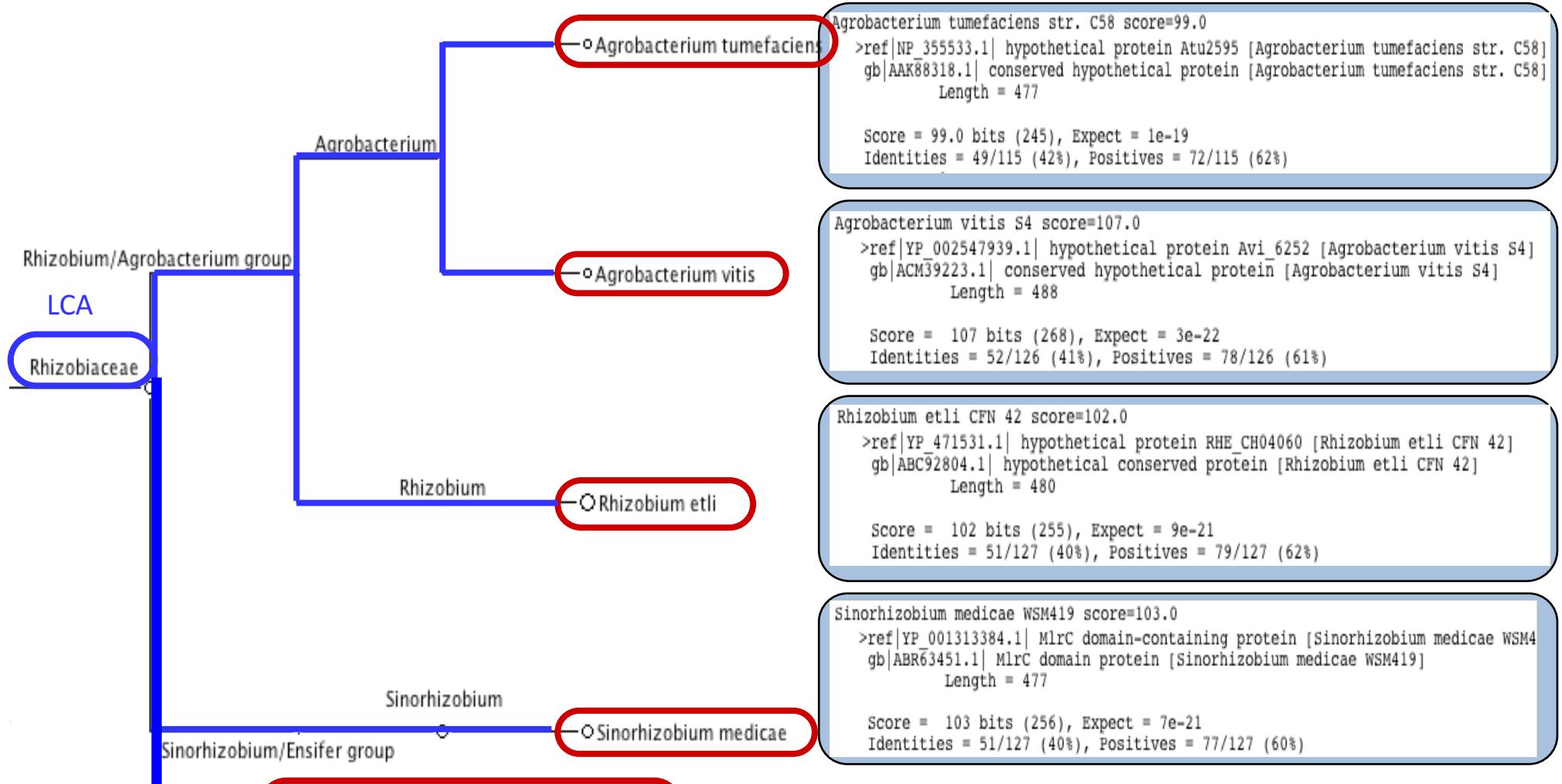


Import Options



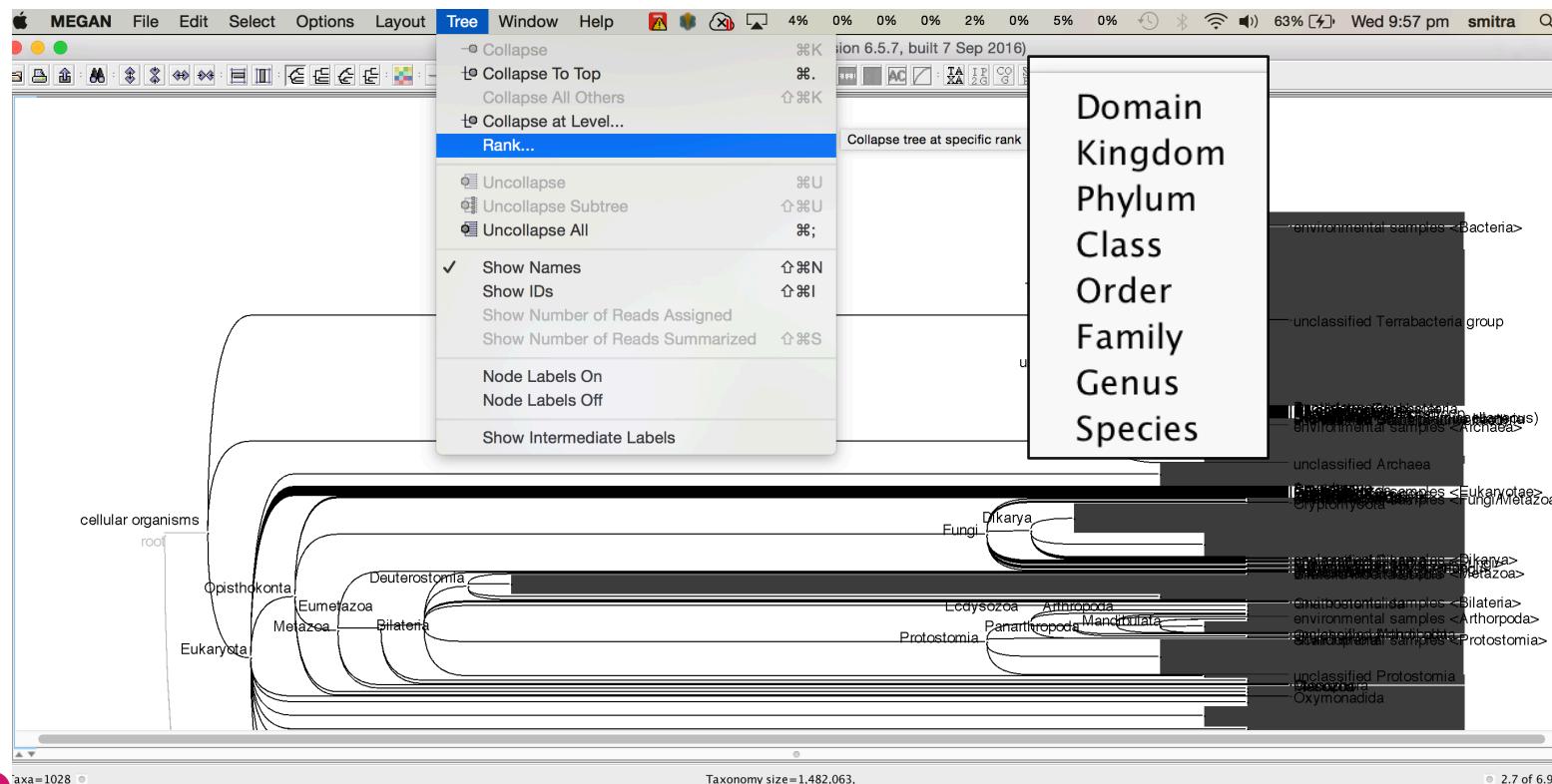
Approach: LCA-Algorithm

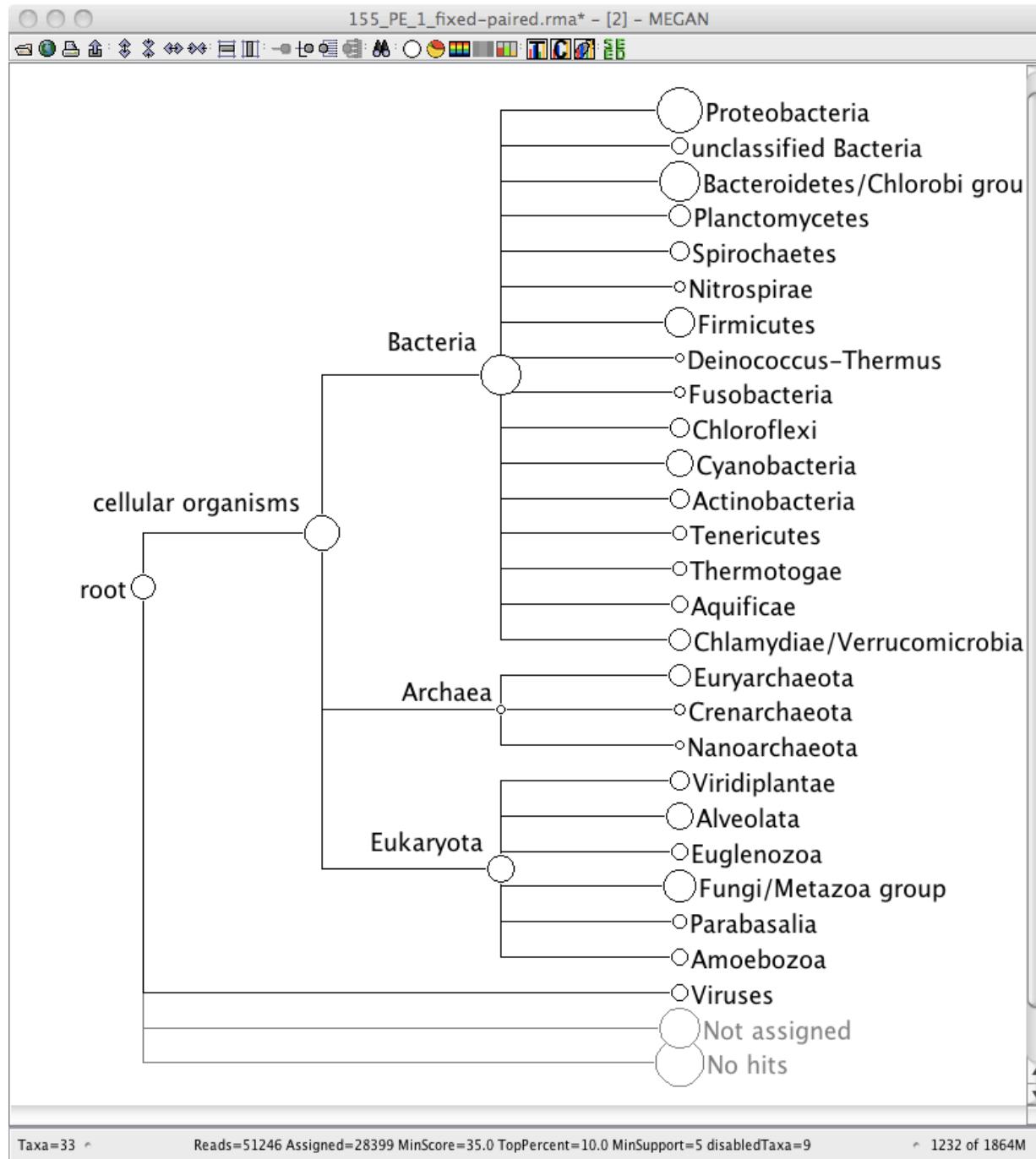
A read will often match more than one database entry



Taxonomic Binning in MEGAN

- MEGAN uses NCBI taxonomy to organize all reads
- NCBI taxonomy tree is currently approx. 1,482,063 nodes and up to 25 ranks deep.





Data Exploration

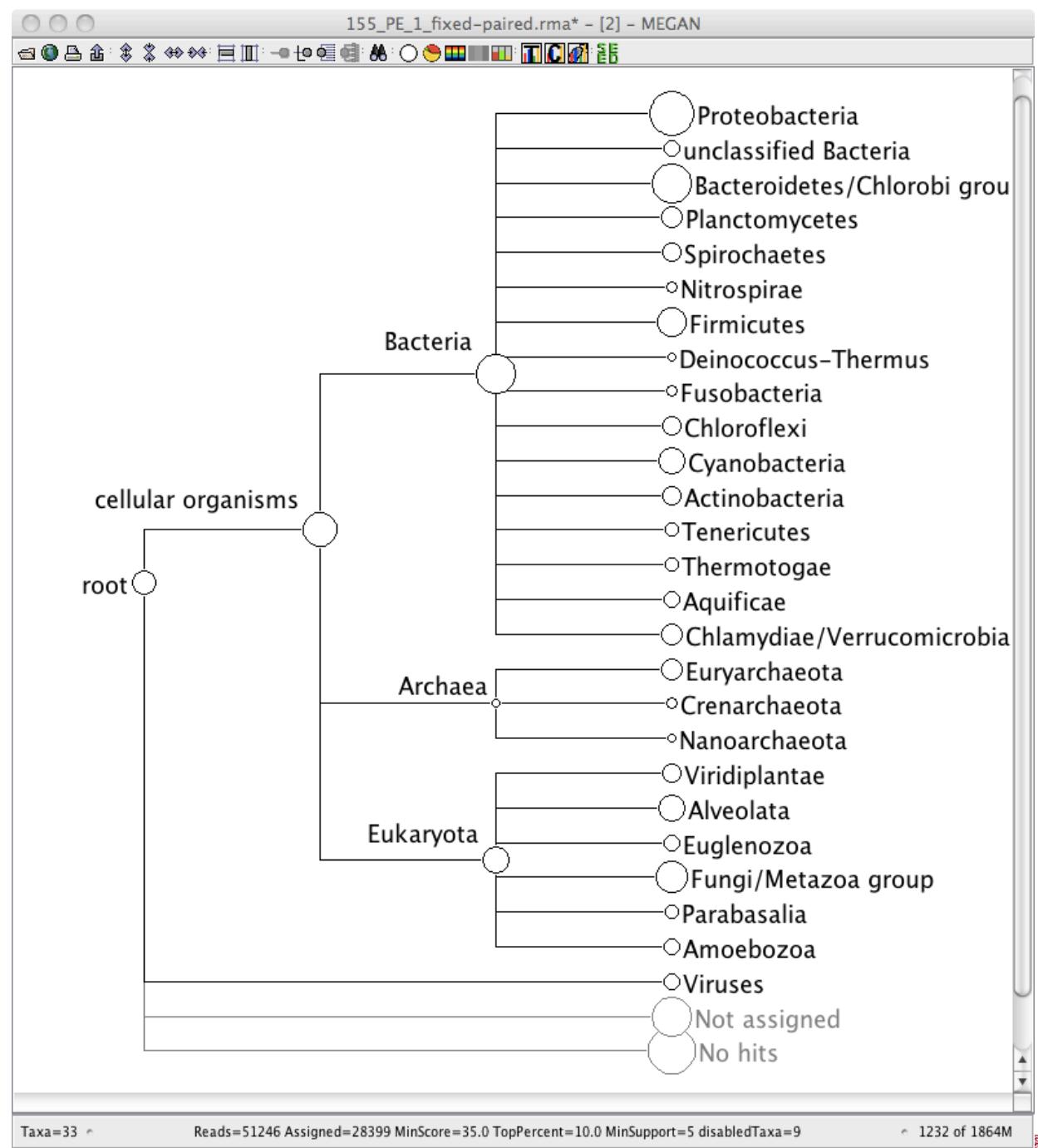
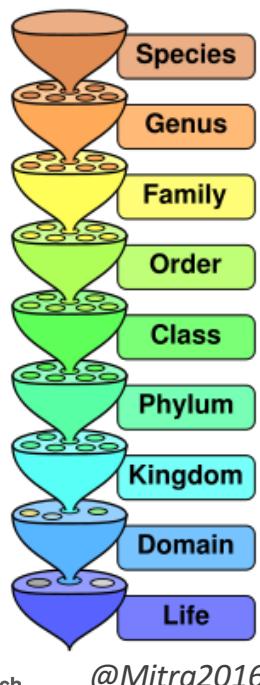
How to analyze a meta genome?

- Organize
- Visualize
- Interact
- Summarize
- Capture
- Compare

Organize and visualize

Taxonomical analysis

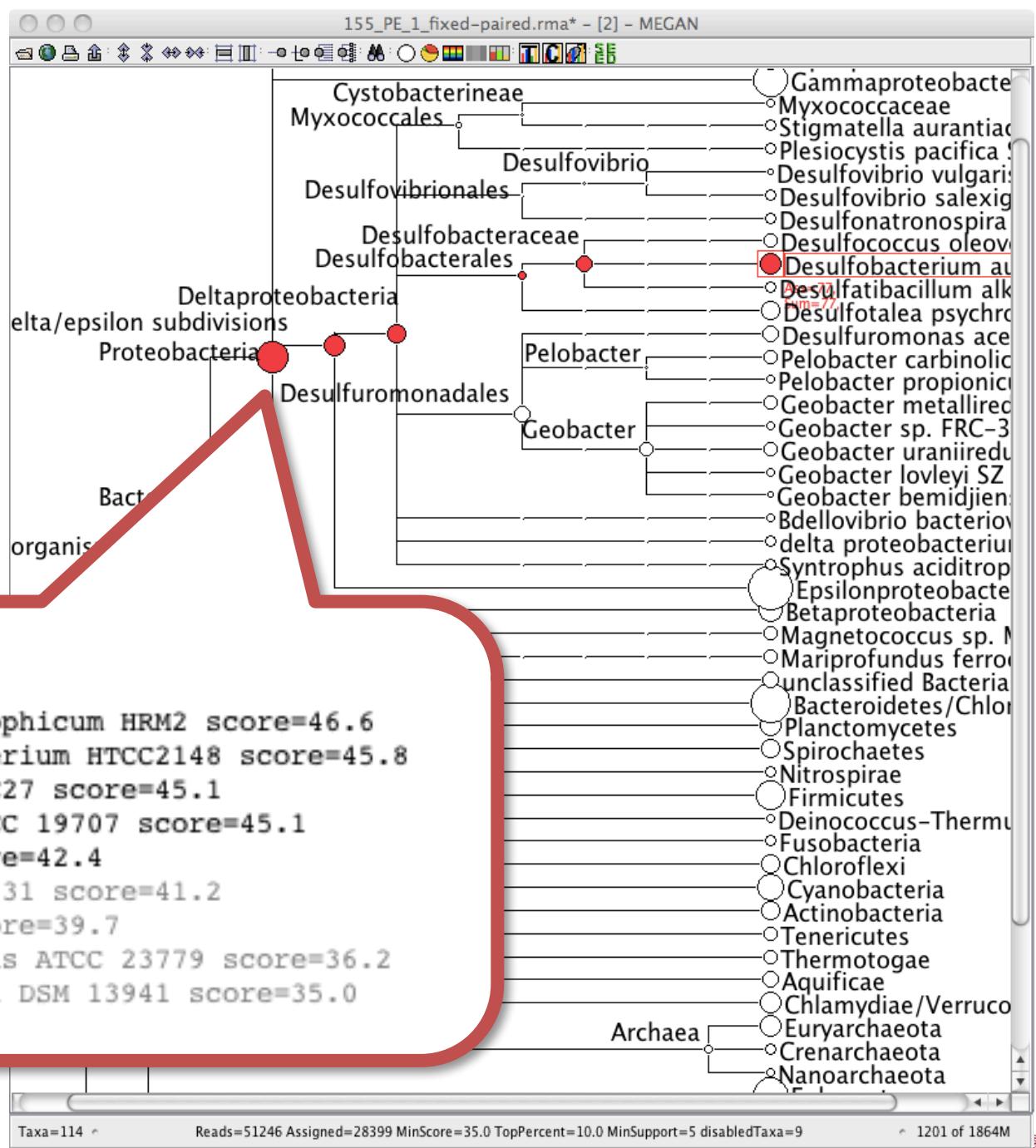
- ✓ Use NCBI taxonomy to structure sequences by evolutionary relatedness of organisms



Organize and visualize Taxonomical analysis

- ✓ Use NCBI taxonomy to structure sequences by evolutionary relatedness of organisms

- ▼ GDEG1CX11G0QZ8.1 [9]
 - ▶ DATA [length=233]
 - ▶ Desulfobacterium autotrophicum HRM2 score=46.6
 - ▶ marine gamma proteobacterium HTCC2148 score=45.8
 - ▶ Nitrosococcus oceani AFC27 score=45.1
 - ▶ Nitrosococcus oceani ATCC 19707 score=45.1
 - ▶ Roseovarius sp. 217 score=42.4
 - ▶ Nitrococcus mobilis Nb-231 score=41.2
 - ▶ Geobacter lovleyi SZ score=39.7
 - ▶ Herpetosiphon aurantiacus ATCC 23779 score=36.2
 - ▶ Roseiflexus castenholzii DSM 13941 score=35.0



Organize and visualize

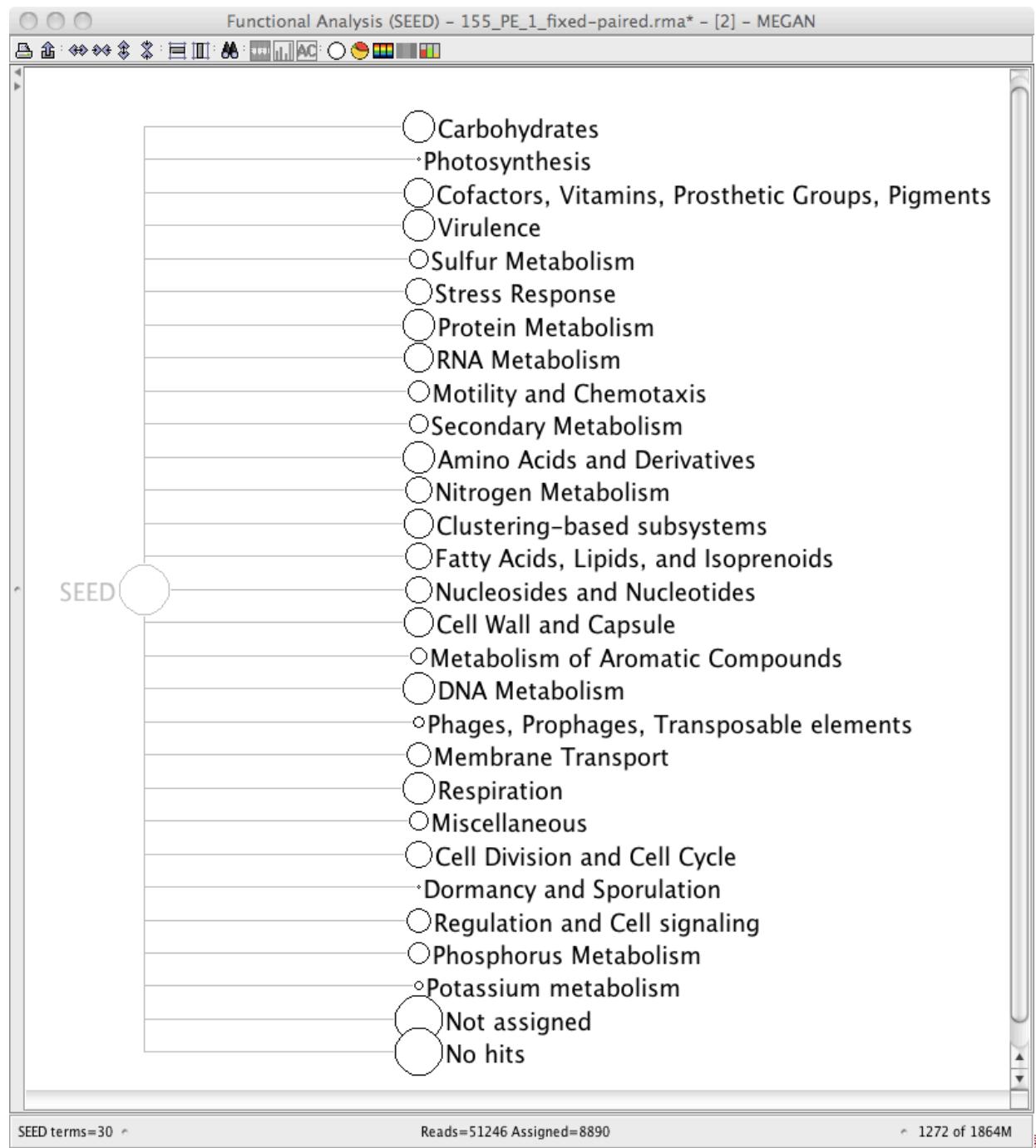
Functional analysis

- ✓ Use SEED classification (like MG-RAST) to structure sequences by subsystems



www.theseed.org

SEED: Overbeek et al.,
Nucleic Acids Res 33 (17), 2005



Organize and visualize

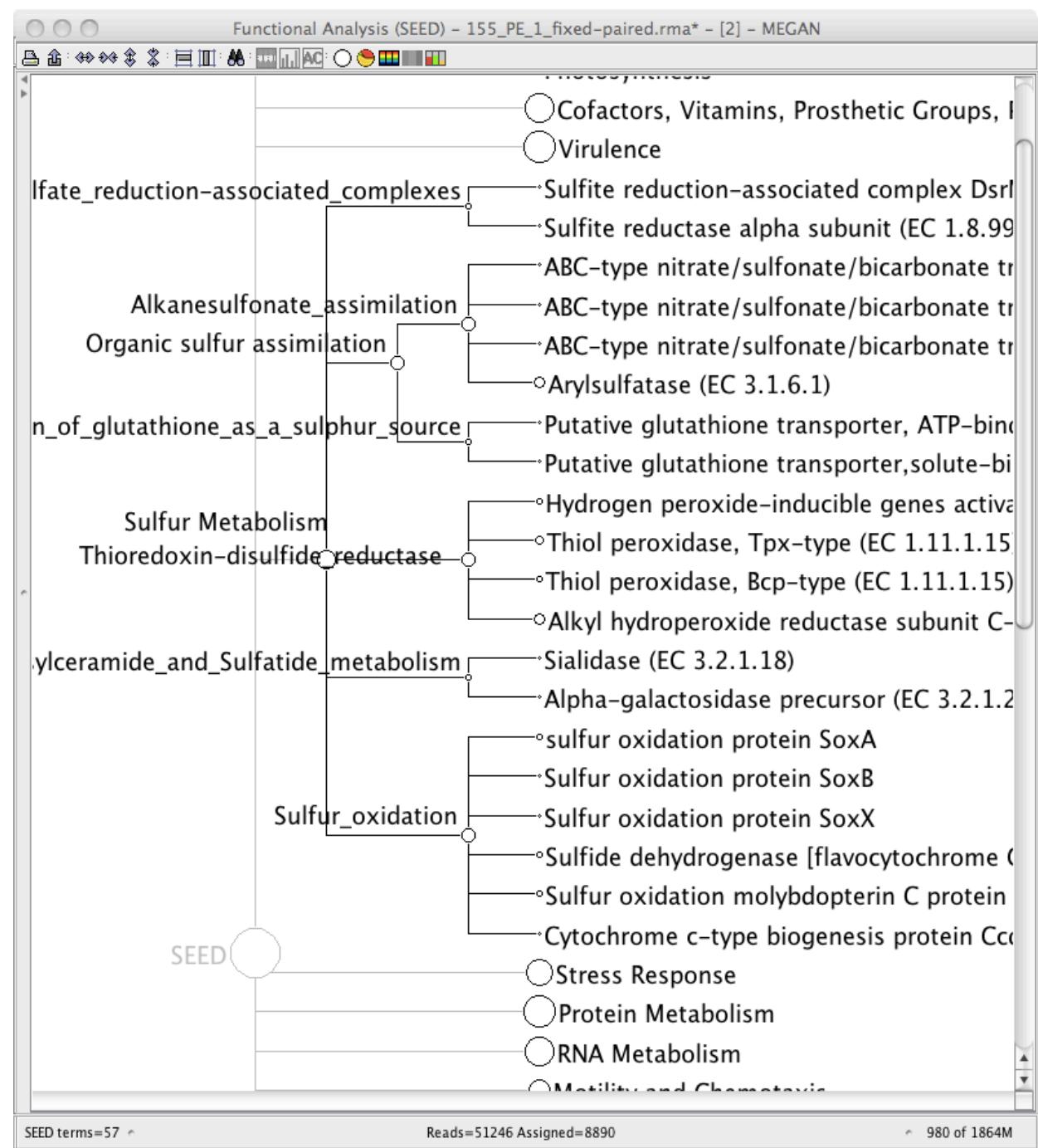
Functional analysis

- ✓ Use SEED classification (like MG-RAST) to structure sequences by subsystems
- ✓ ... and by functional roles



www.theseed.org

SEED: Overbeek et al.,
Nucleic Acids Res 33 (17), 2005



Organize and visualize

CARBON FIXATION IN PHOTOSYNTHETIC ORGANISMS

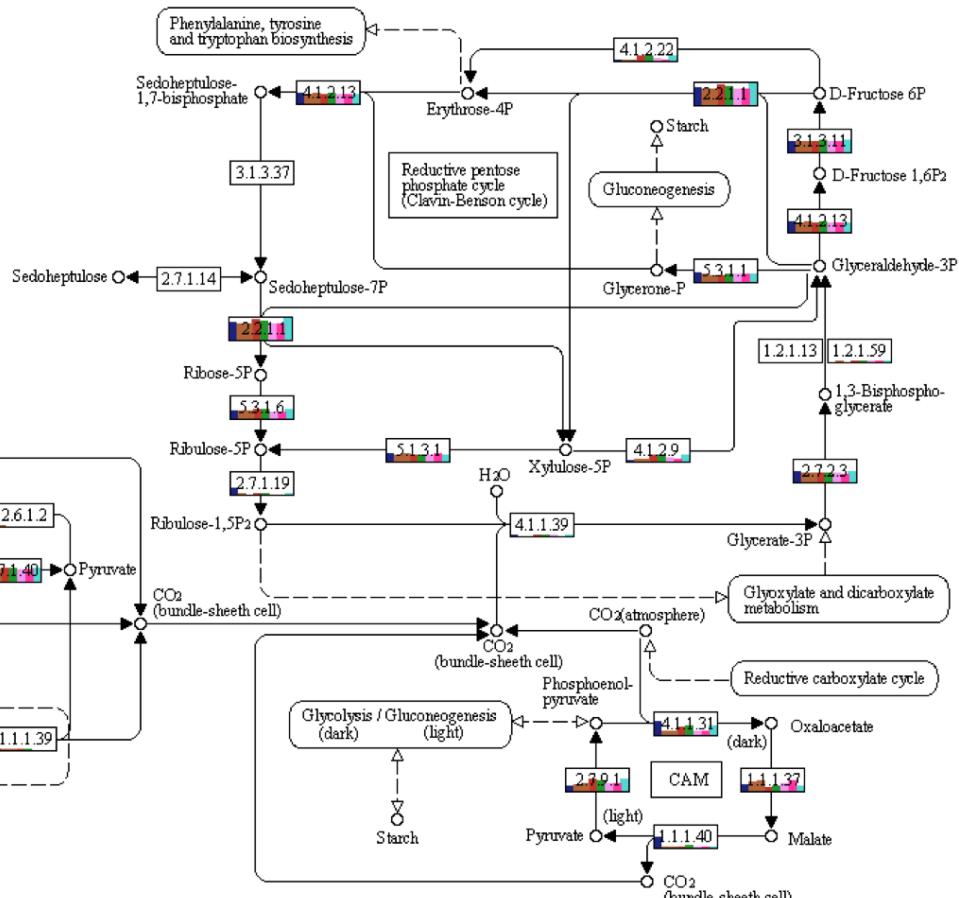
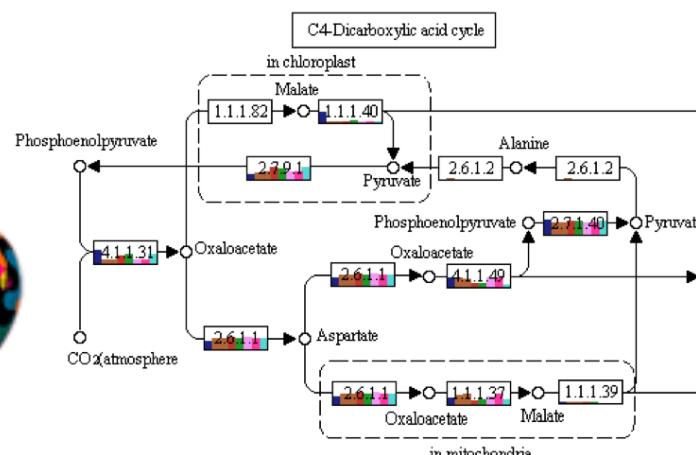
Functional analysis

- ✓ Use KEGG pathways to structure sequences by their presence in pathways



<http://www.genome.jp>

KEGG: Kanehisa et al,
Nucleic Acids Res. 38, D355-D360 (2010)



MEGAN KEGG VIEWER

Various Chart Views

Clostrid
Paenibacillaceae

186_anodeb

Clostrid
Paenibacillaceae
Bacill

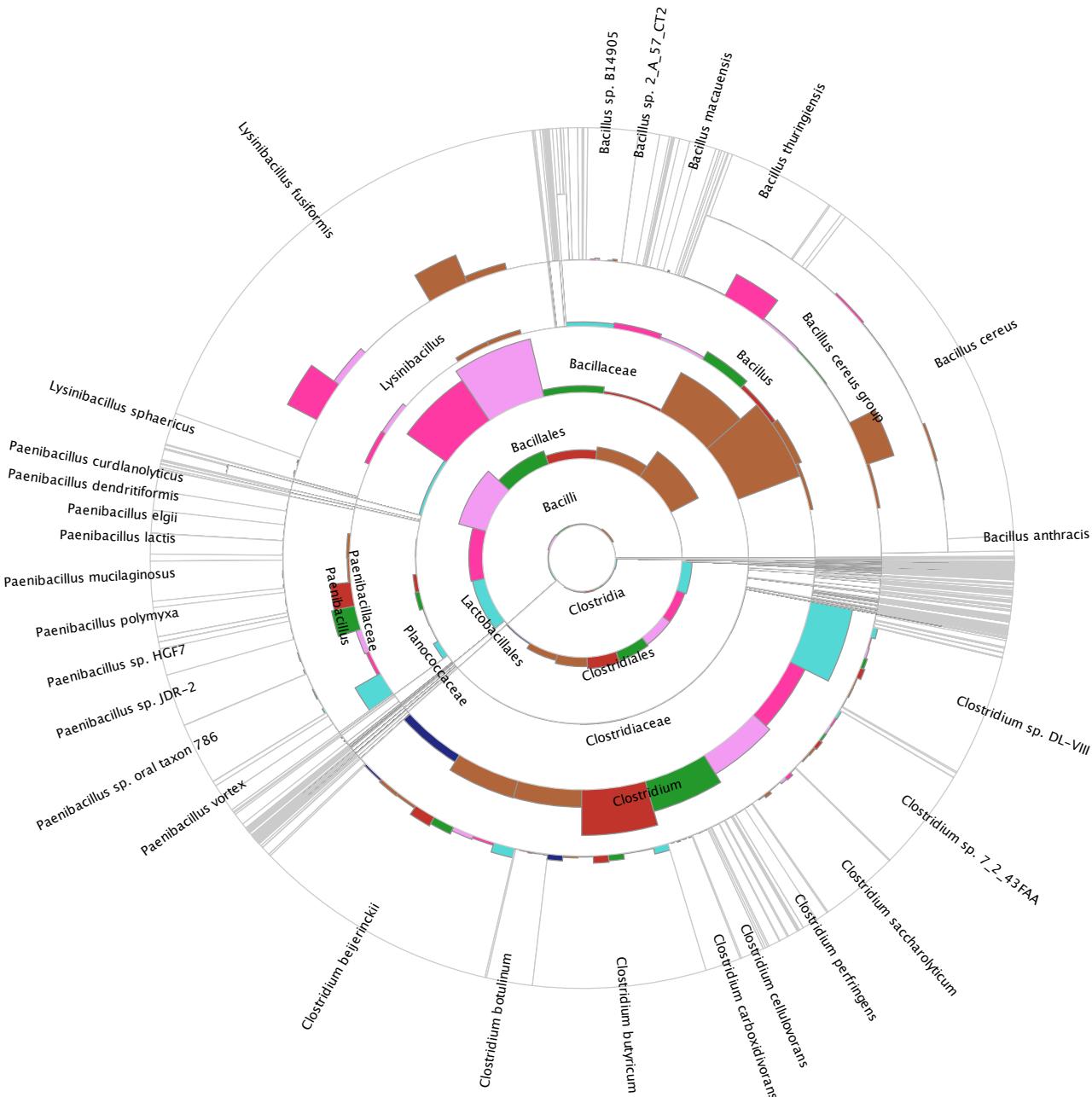
183_anodeb

Enterobacte
Clostridiaceae
Bacill

533_anoly

Clostridiac
Paenibacillaceae
Enterobacteriaceae

538_anodebi



illaceae
Clostridiaceae

189_anolyte_new

bacillaceae
Clostridiaceae

2_anolyte_A1B1

ridiaceae
cillaceae
Enterobacteriaceae

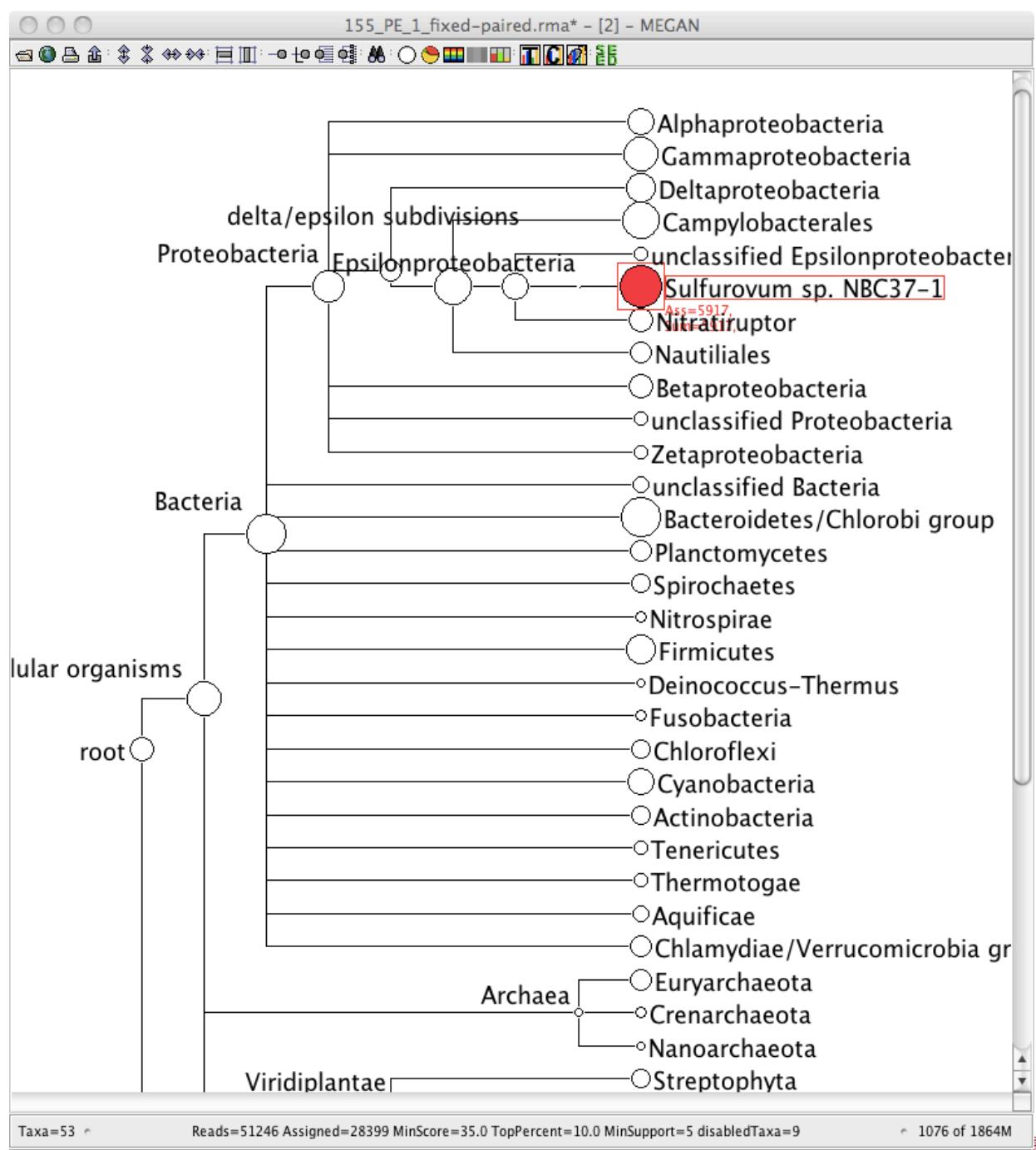
7_anolyte_B6B1

cillaceae
obacteriaceae
Clostridiaceae

anodebiofilm_B6A

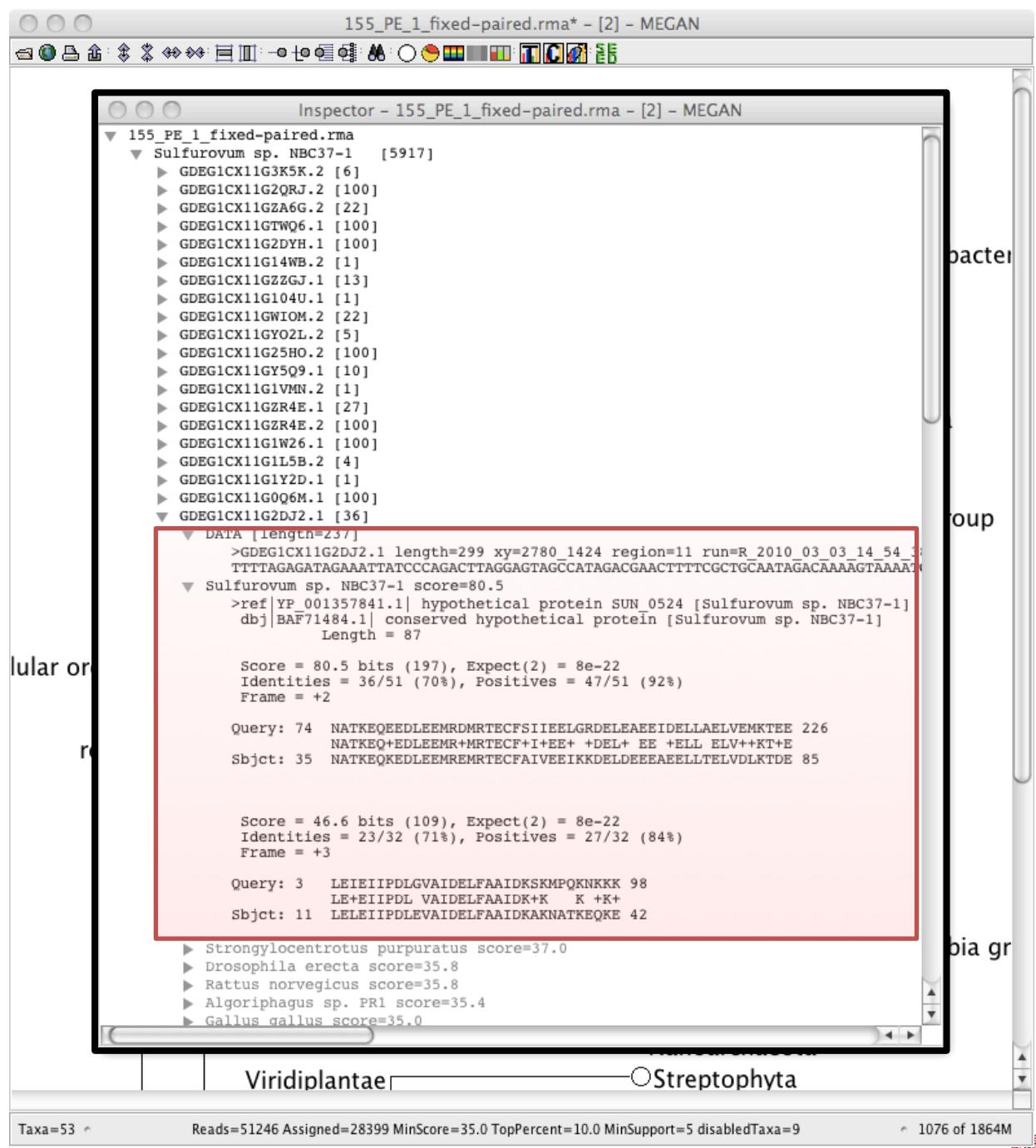
Interact and summarize

- ✓ Search for nodes of interest



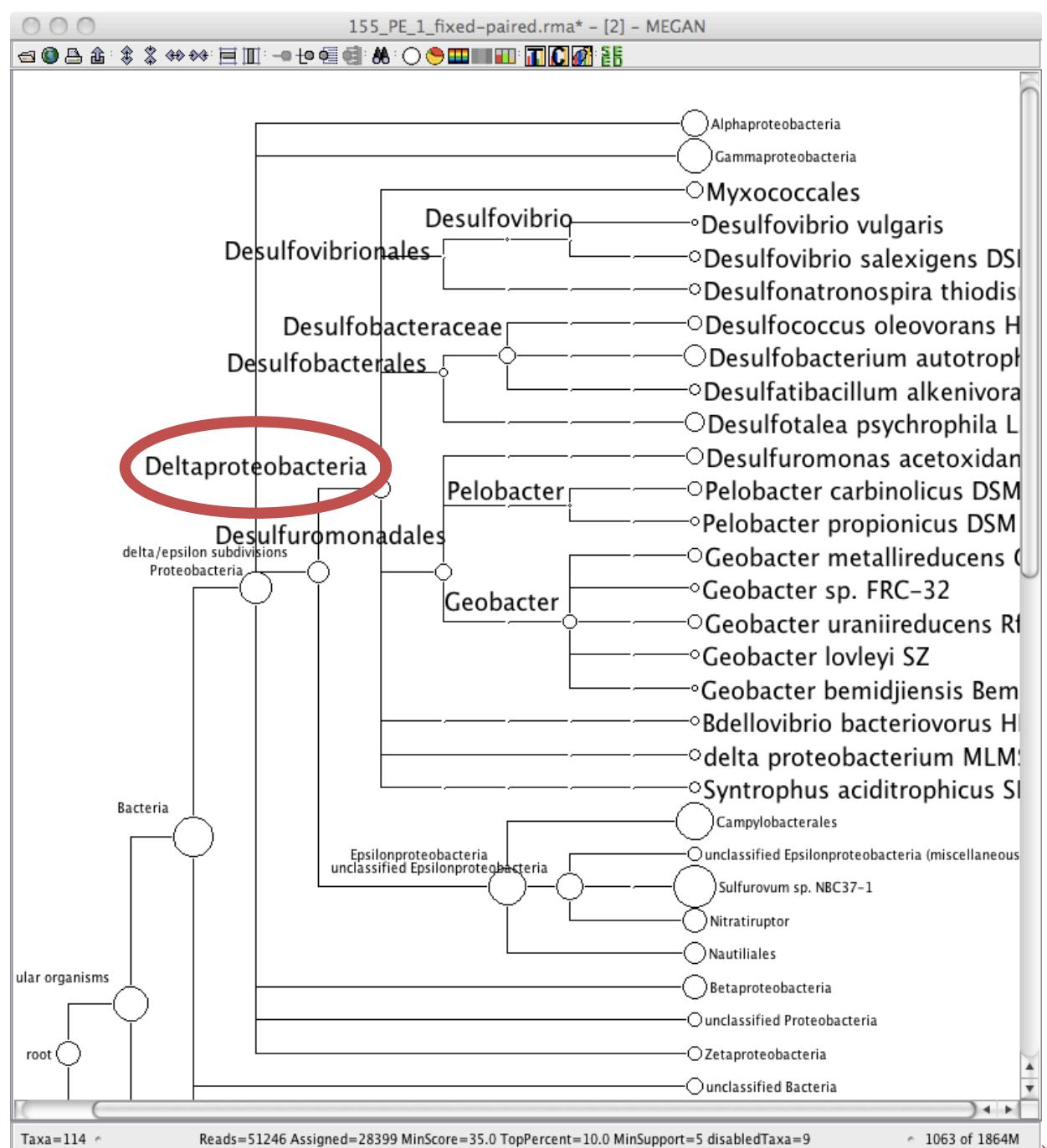
Interact and summarize

- ✓ Search for nodes of interest
- ✓ Inspect sequences assigned to a node



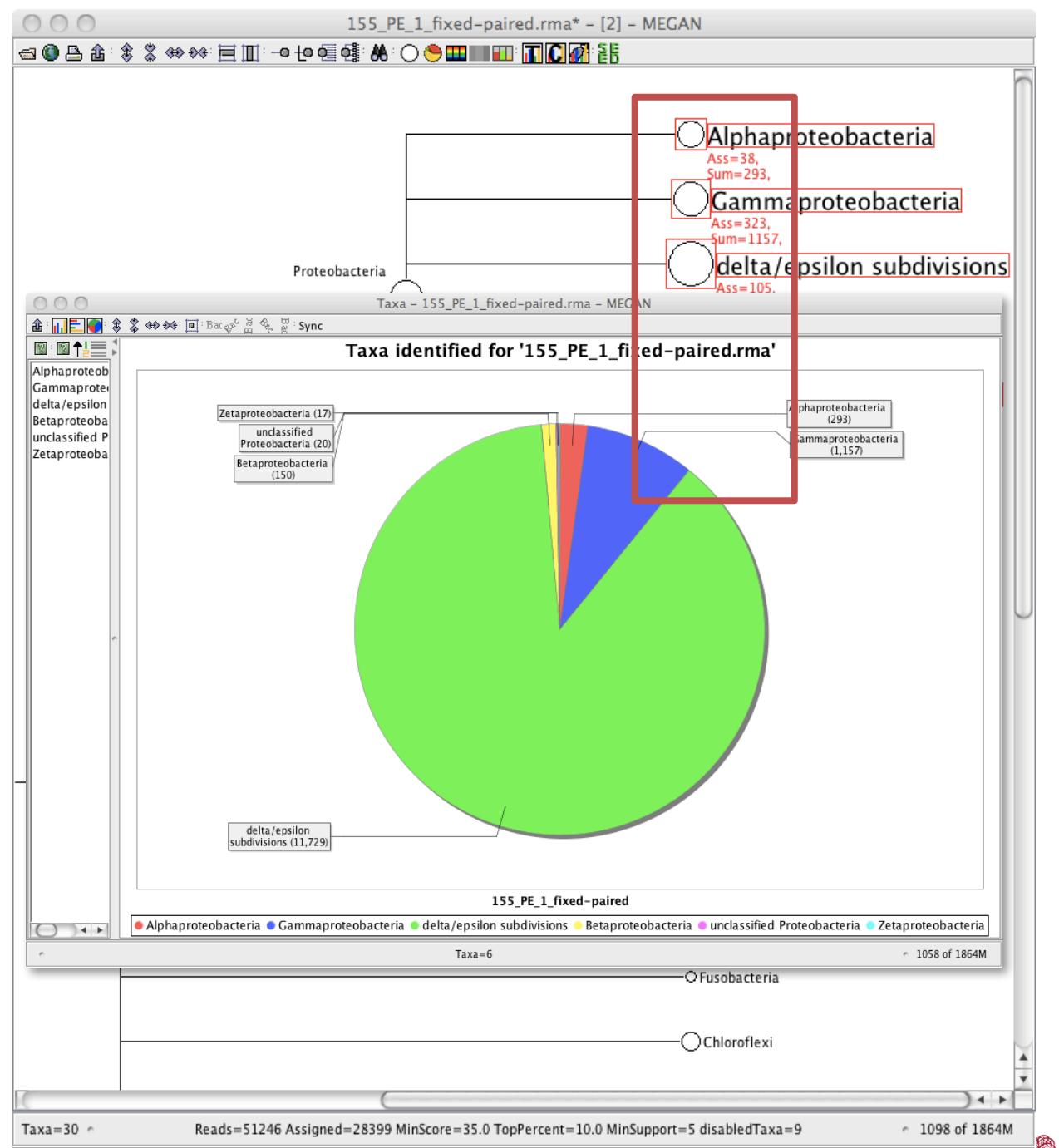
Interact and summarize

- ✓ Search for nodes of interest
- ✓ Inspect sequences assigned to a node
- ✓ Collapse and un-collapse parts of the tree



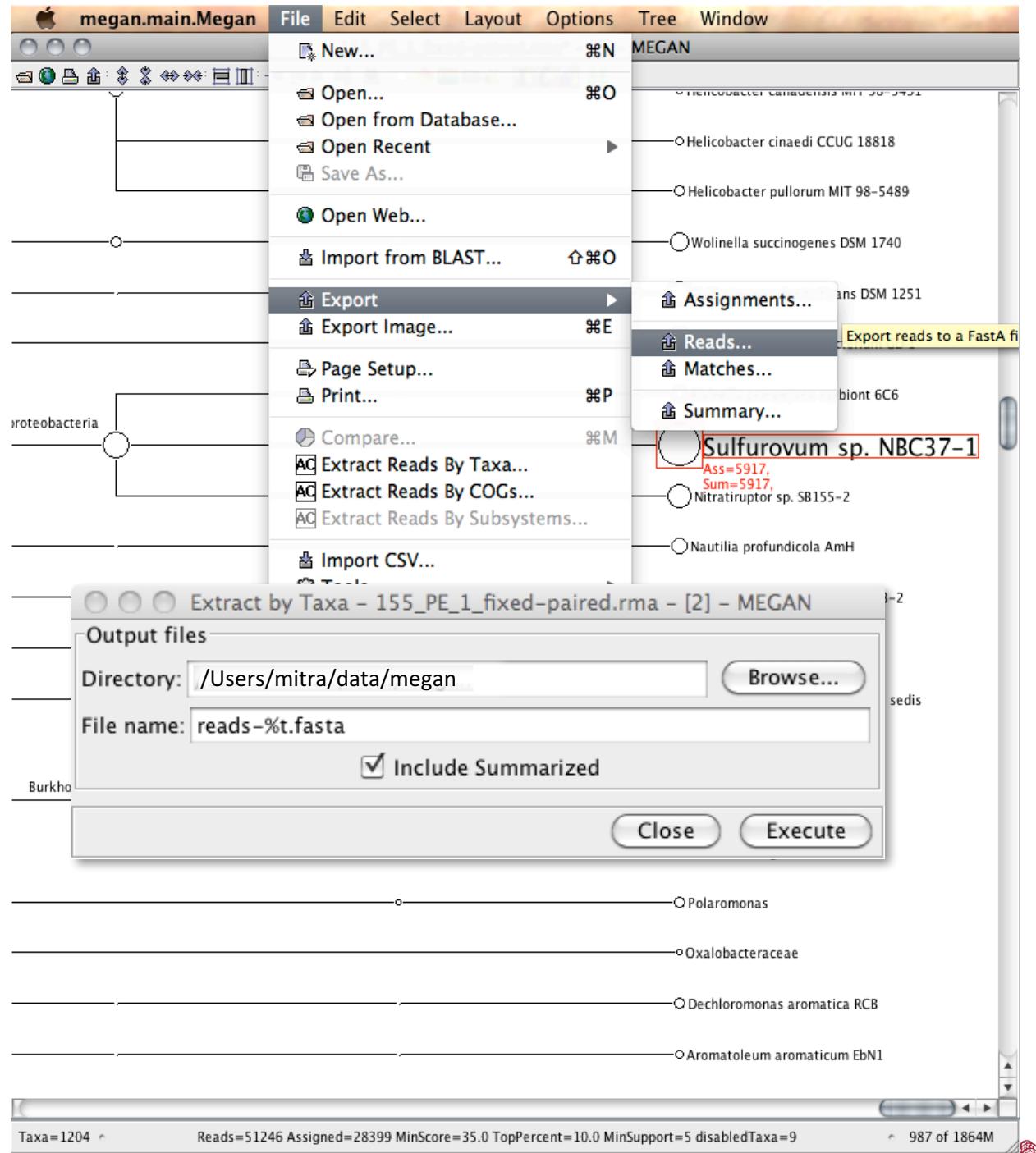
Interact and summarize

- ✓ Search for nodes of interest
- ✓ Inspect sequences assigned to a node
- ✓ Collapse and un-collapse parts of the tree
- ✓ Create charts



Capture

- ✓ Capture all sequences (and/or their matches) assigned to selected nodes
- ✓ Taxonomy
- ✓ SEED
- ✓ KEGG



Comparative Metagenomics

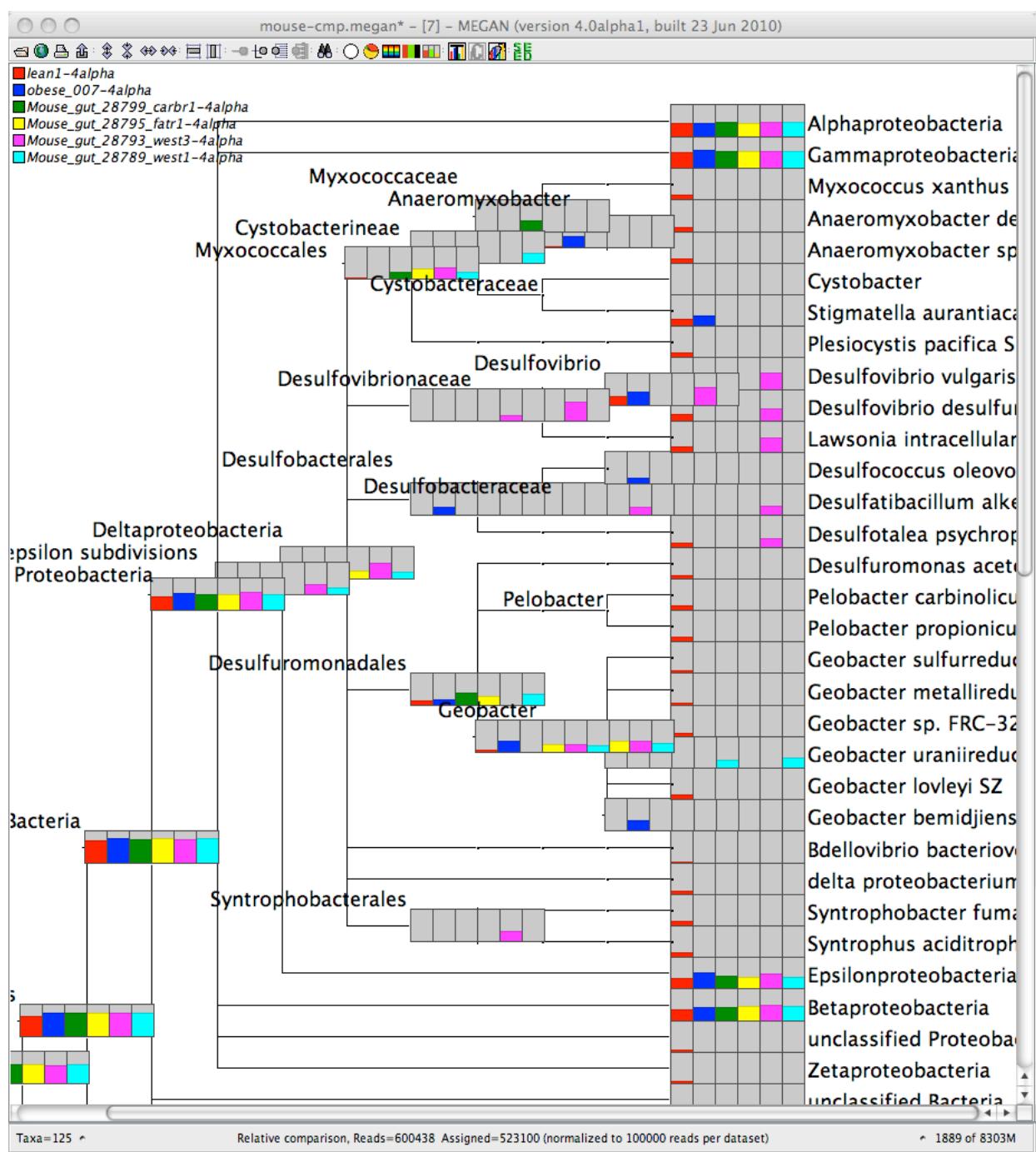
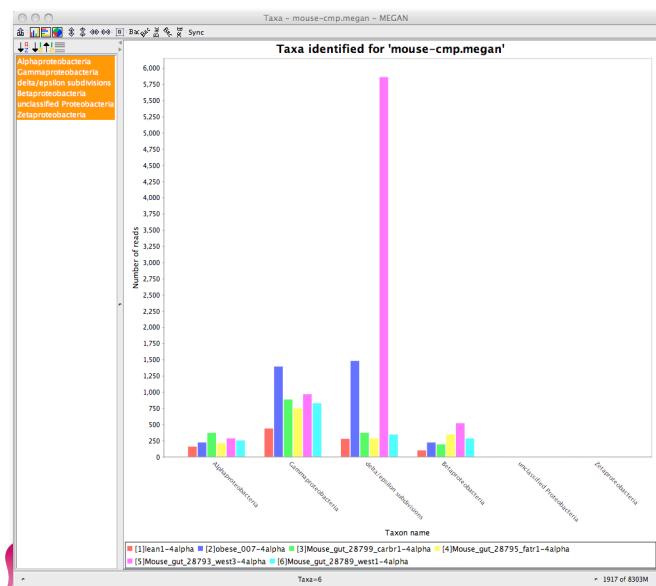
- Why would one want to?
 - Structure of microbial communities
 - Time series studies
 - Based on geography
 - Clinical research to correlate microbes with disease
 - Comparing individuals
 - Different time point
 - Different drug influence
 -
- Methods for the comparison is needed



Visual Comparison

Display multiple datasets simultaneously

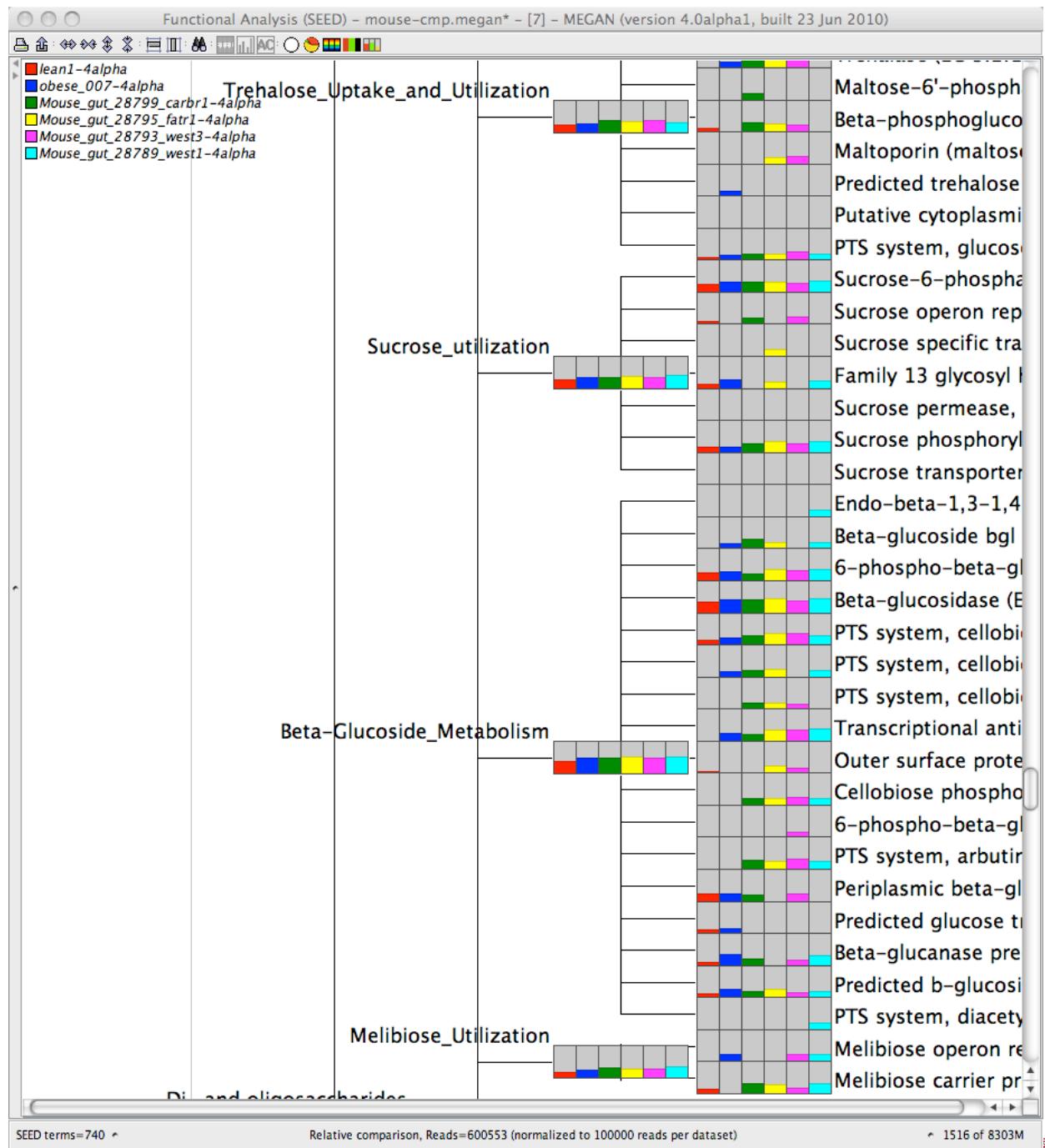
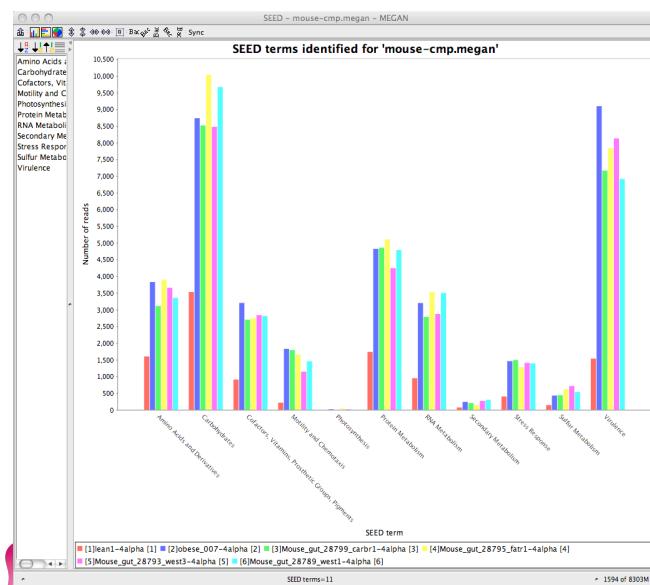
- ✓ Taxonomical comparison
- ✓ Interact
- ✓ ... and summarize



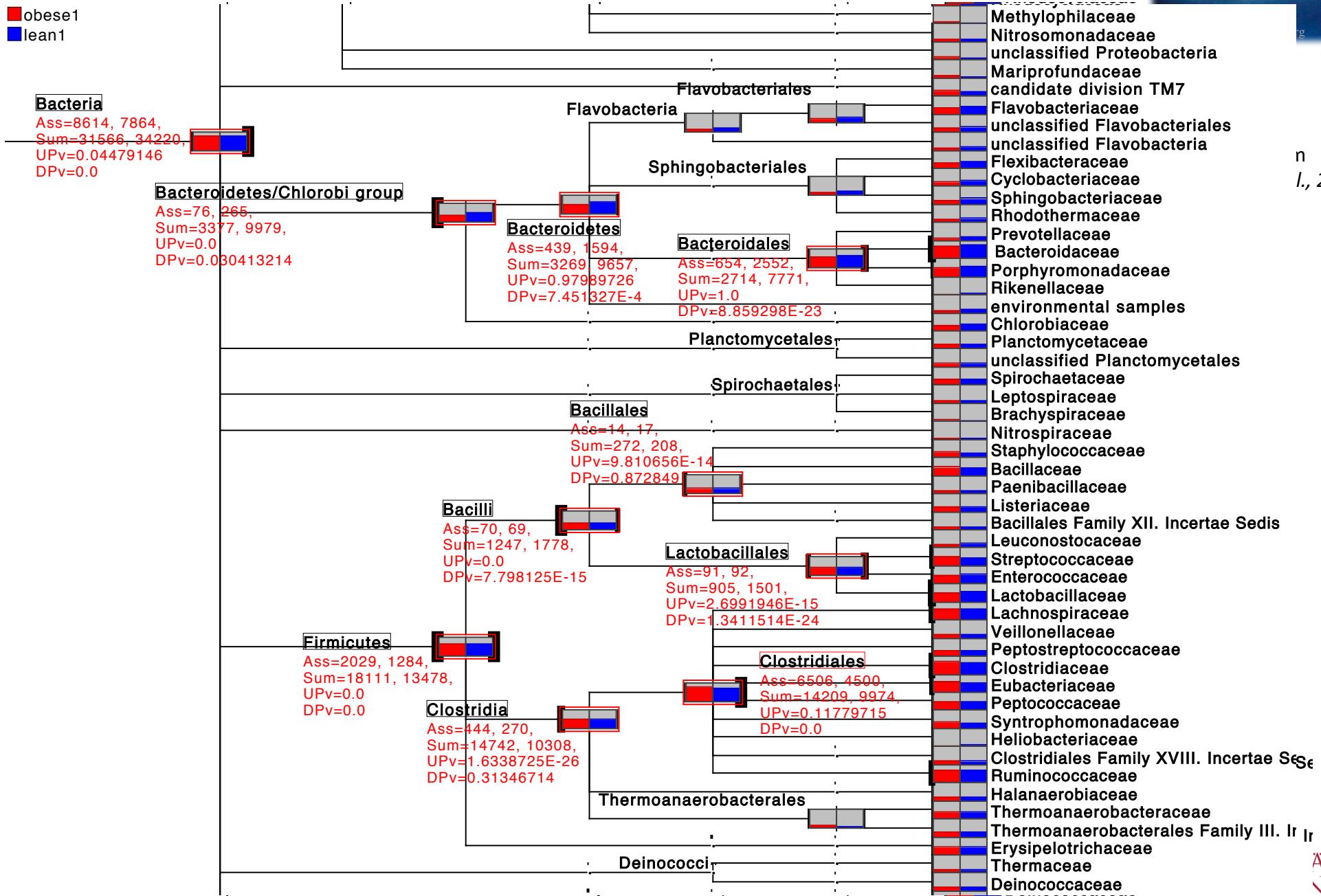
Visual Comparison

Display multiple datasets simultaneously

- ✓ Taxonomical comparison
- ✓ Interact
- ✓ ... and summarize
- ✓ Functional comparison



Pair-wise Comparison



n
l., 2006

AT



Multiple Comparison

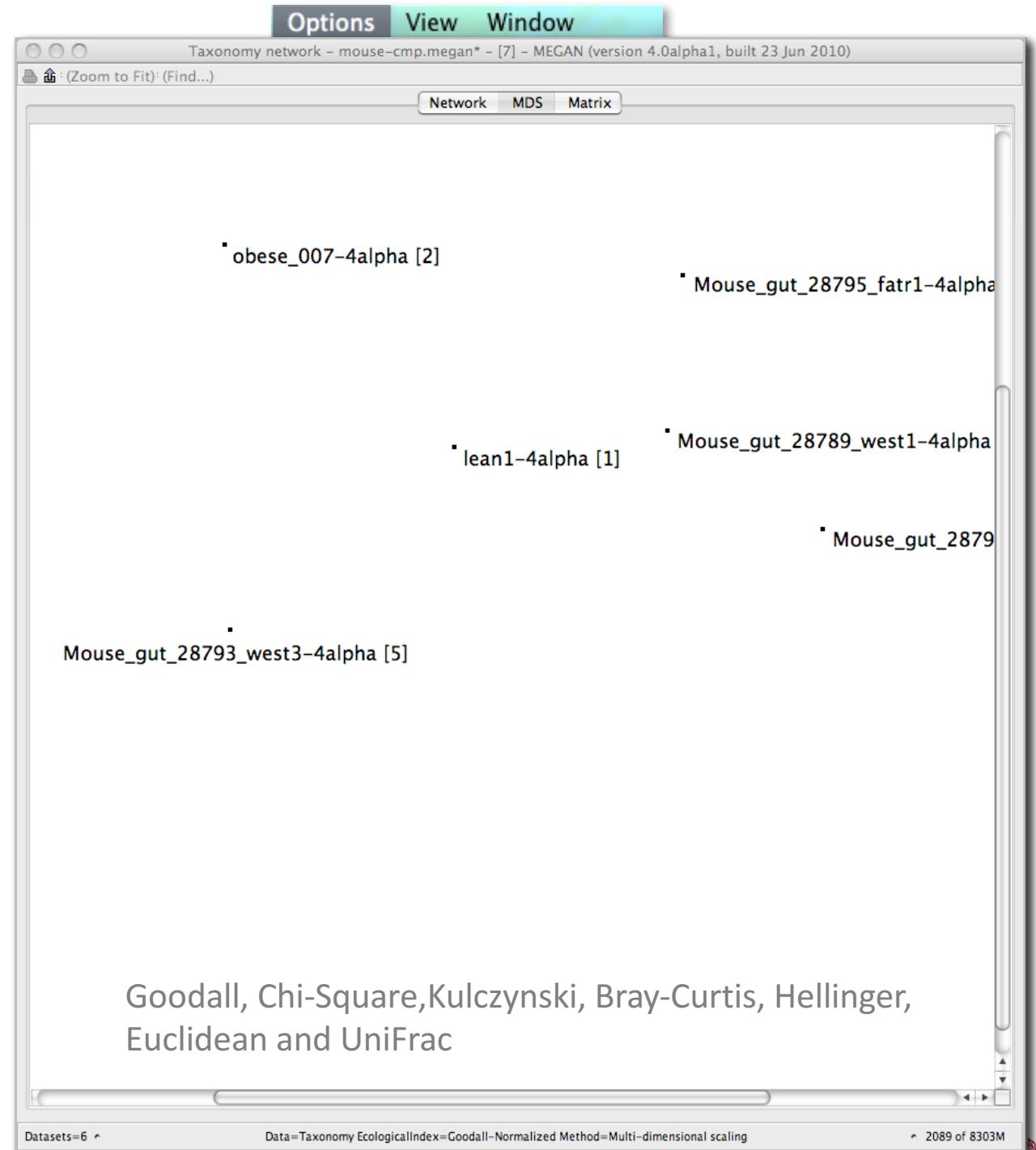
High-level comparison:

- ✓ Select taxa
- ✓ Compute ecological indices (distances)
- ✓ Represent distances using neighbor-net
- ✓ ... or MDS

Mitra, Gilbert, Field and Huson,
ISME J, 2010

Neighbor-net:
Bryant and Moulton, 2003

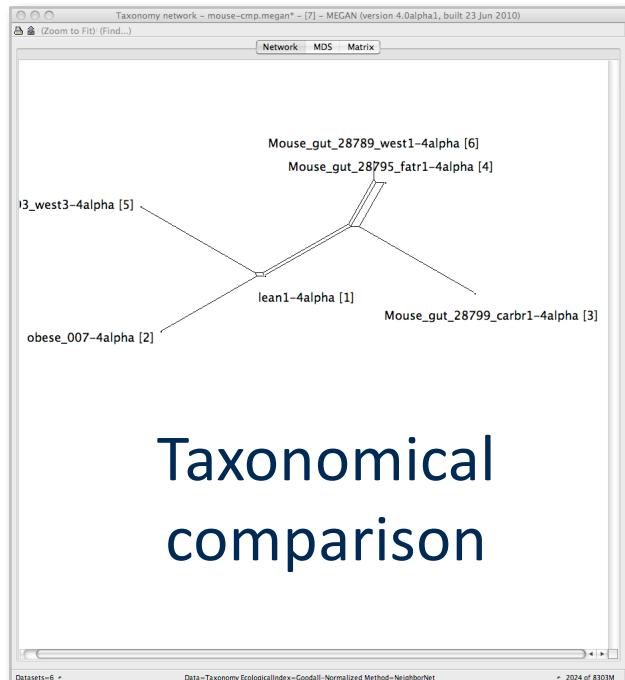
Goodall, Chi-Square, Kulczynski, Bray-Curtis, Hellinger,
Euclidean and UniFrac



Multiple Comparison

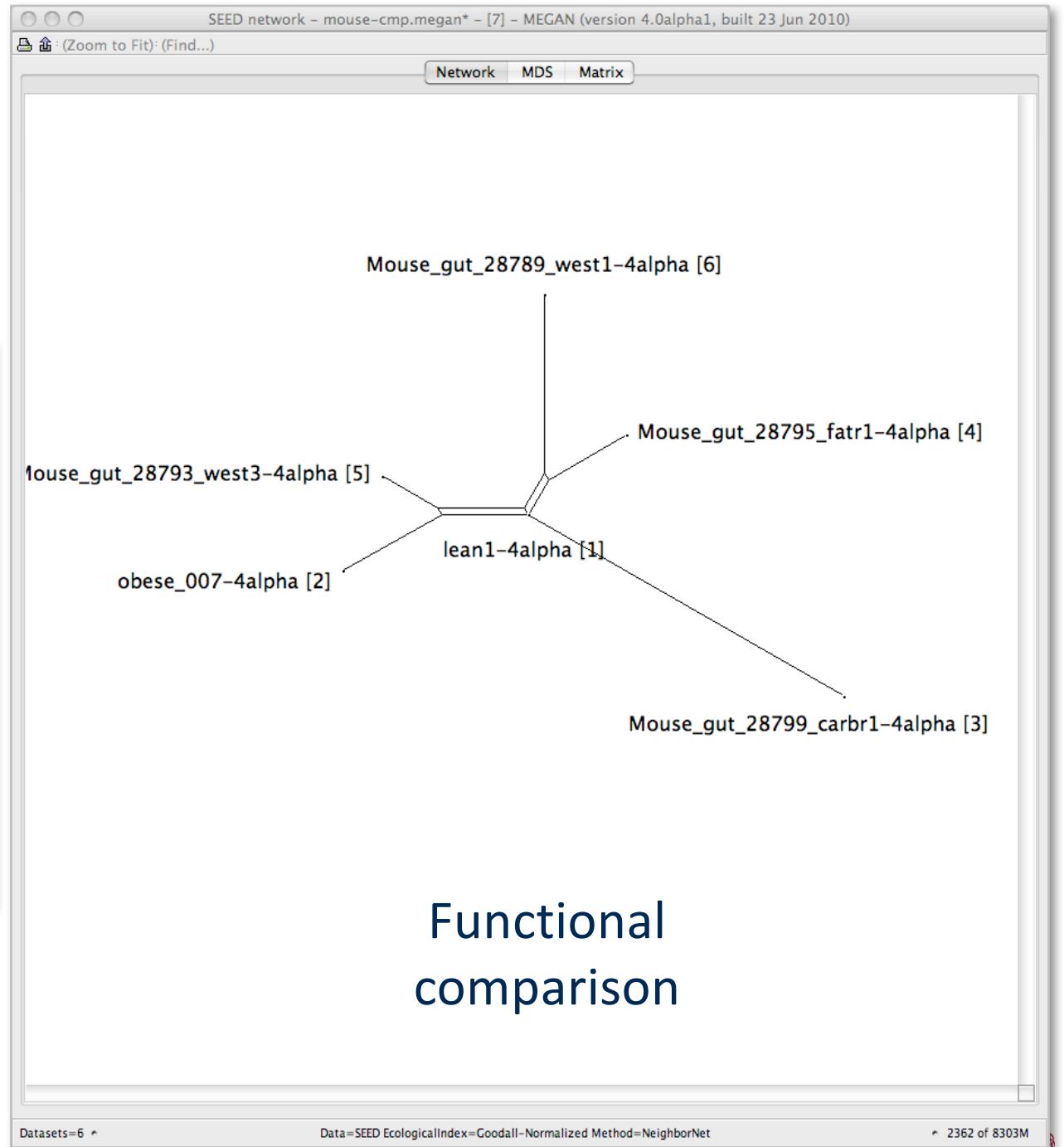
High-level comparison:

- ✓ Select taxa



Taxonomical
comparison

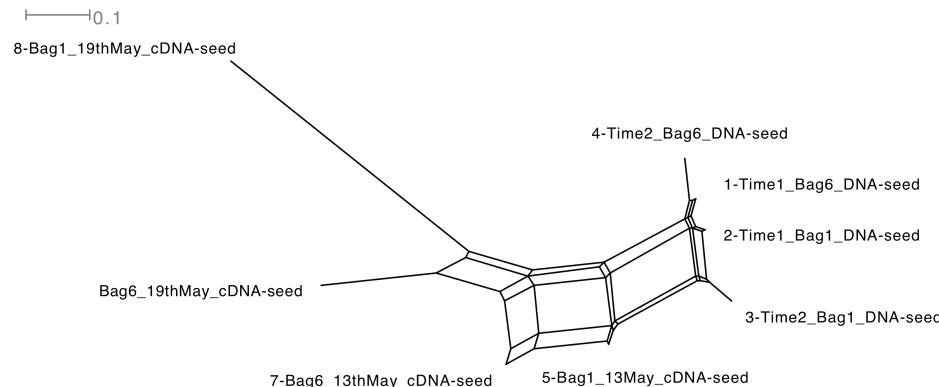
- ✓ Select functions...



Functional
comparison

Network Visualization

- Metagenomes do not evolve along a tree
- Numerous environmental factors may result in incompatible signals
- Neighbor-net method is used to compute an unrooted network between multiple datasets

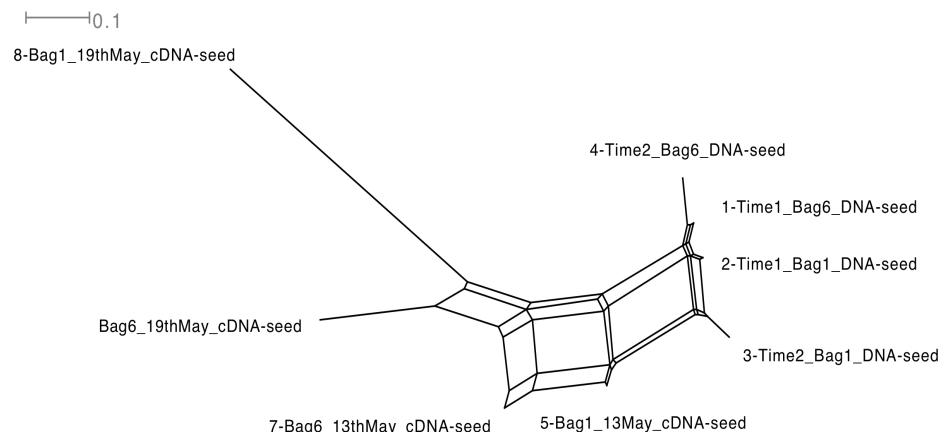


Mitra et al., 2010

Computational Approach

- Data matrix obtained from taxonomic or functional profile
- Distances obtained from pair combinations using different ecological indices

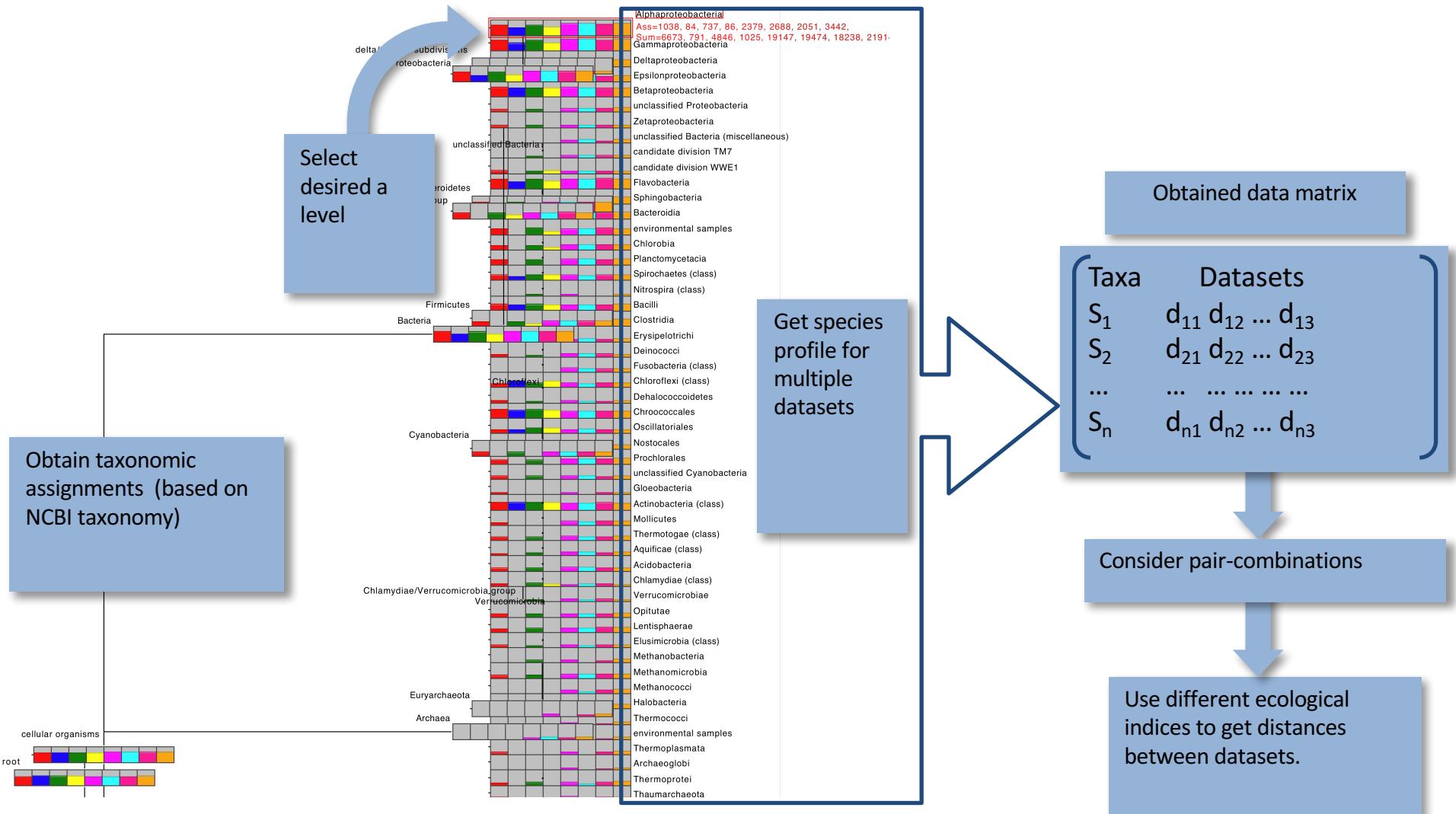
$$\mathcal{D} = \begin{pmatrix} d_{11} & d_{12} & \dots & \dots & \dots & d_{1n} \\ \vdots & d_{22} & & & & \vdots \\ \vdots & \vdots & \ddots & & & \vdots \\ \vdots & \vdots & & \ddots & & \vdots \\ \vdots & \vdots & & & \ddots & \vdots \\ d_{n1} & d_{n2} & \dots & \dots & \dots & d_{nn} \end{pmatrix}.$$



Mitra et al., 2010



Computational Approach



Performance

- Euclidean distance
- Kulczynski distance
- Bray-Curtis distance
- Chi-squared distance
- Hellinger distance
- Goodall's index

} Dominated by the large values

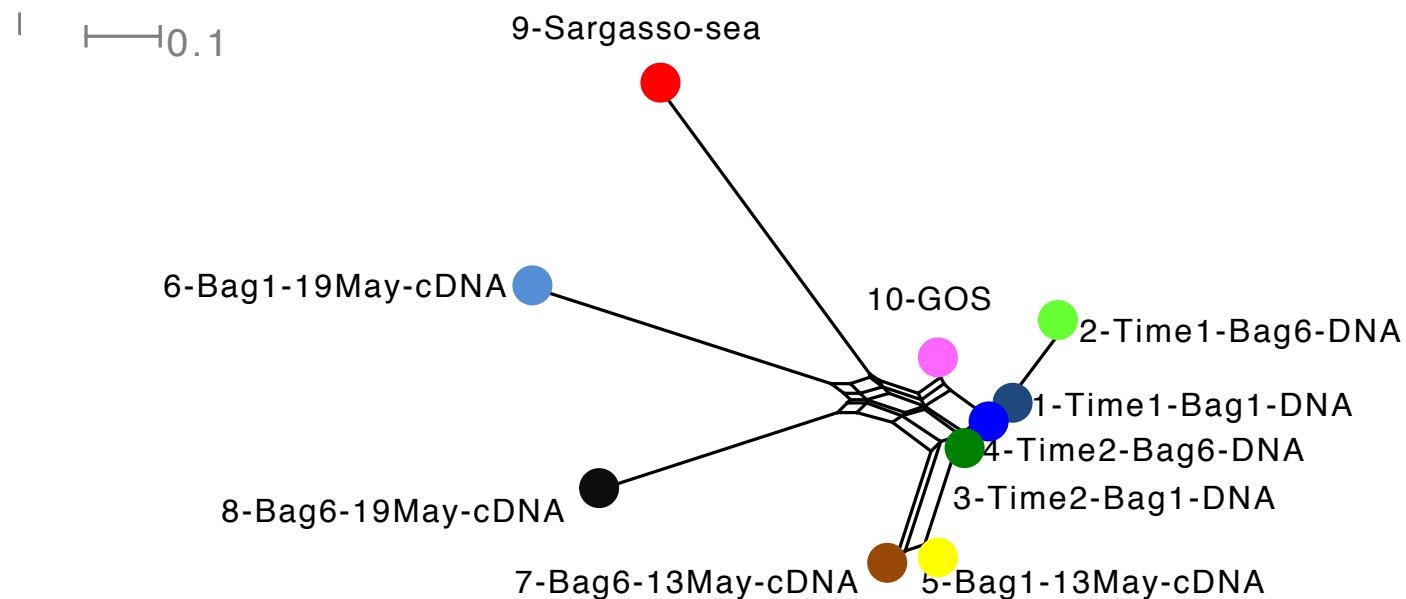
} Based on differences in the proportions of values

} Non-parametric measure for multivariate datasets (effective for the study of rare-taxa)

Robustness

Comparison of 10 marine datasets

Similarity analysis at different levels of taxonomy using ecological indices



Top level

Phylum

Class

Order

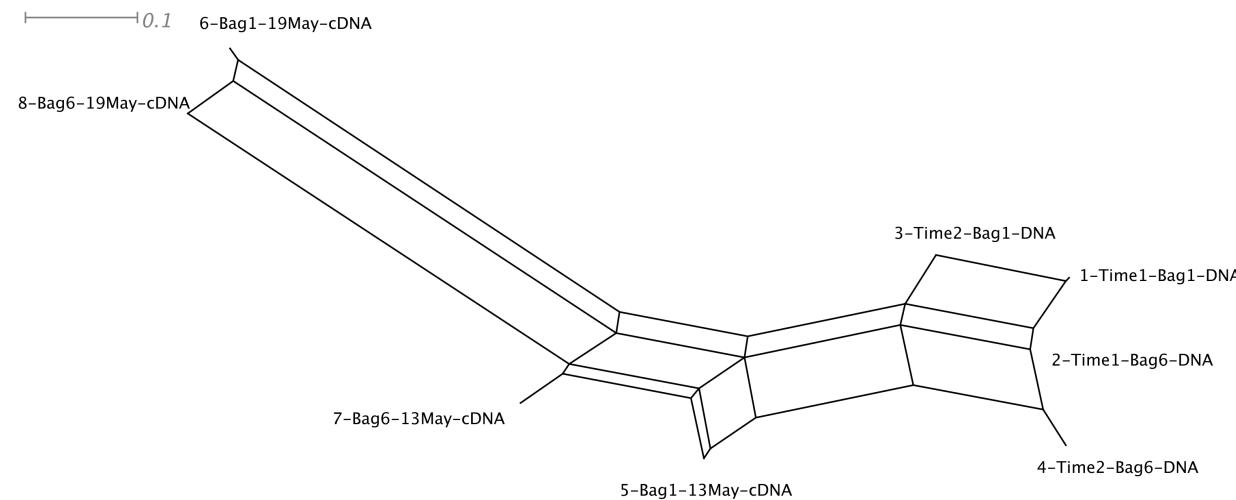
Family

Genus

Species

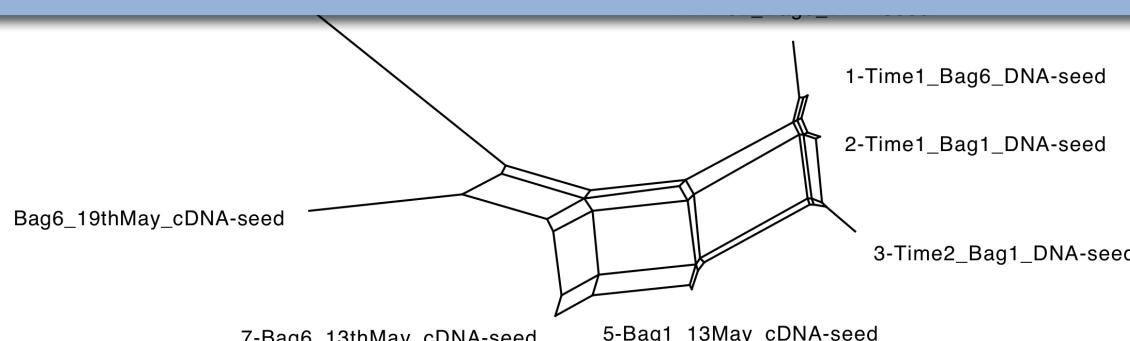
Robustness

Comparison of eight marine metagenome



Network obtained from
taxonomic profile

Clusters are quite similar for both taxonomic and functional content...

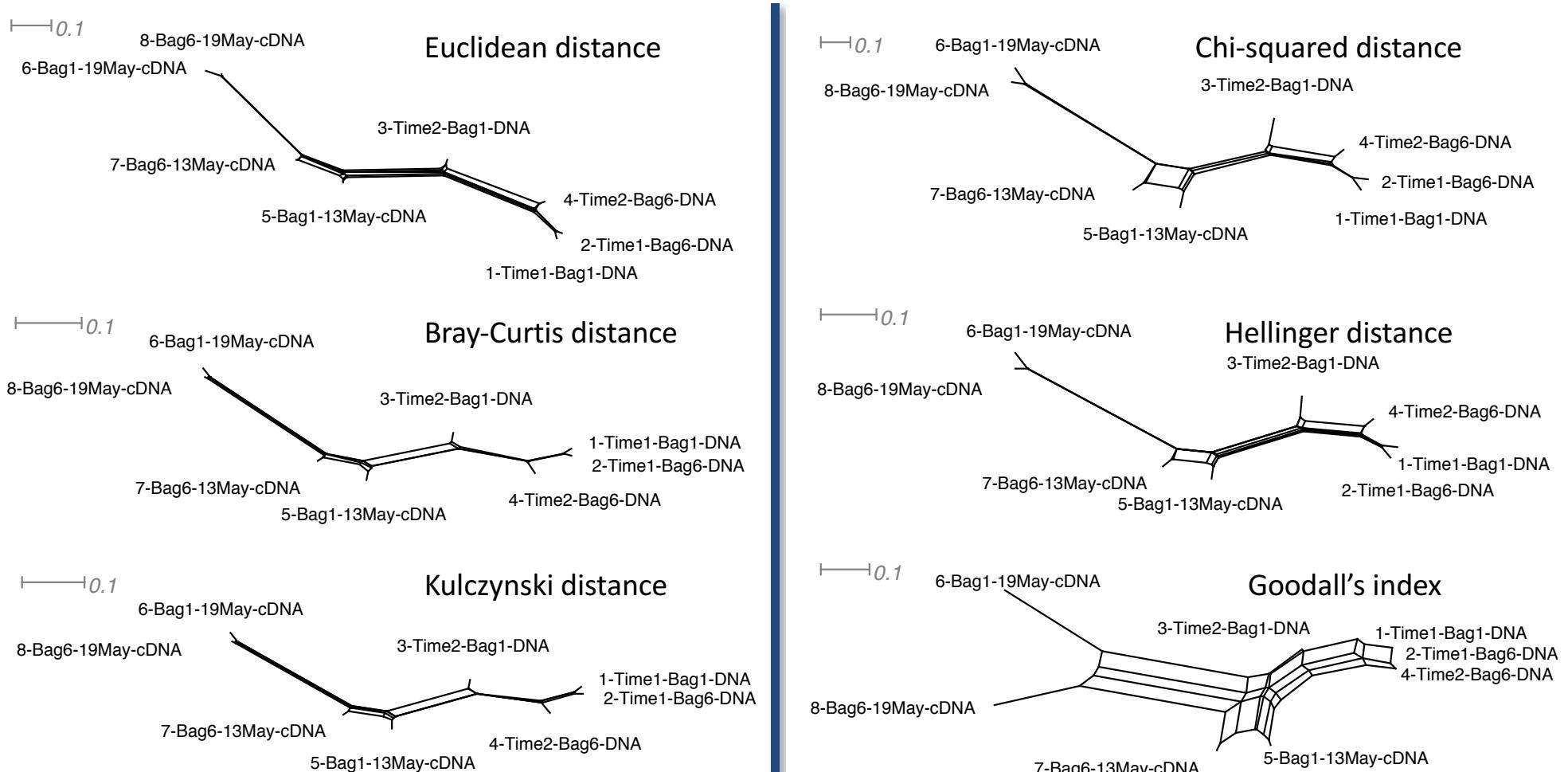


Network obtained from
functional profile



Performance

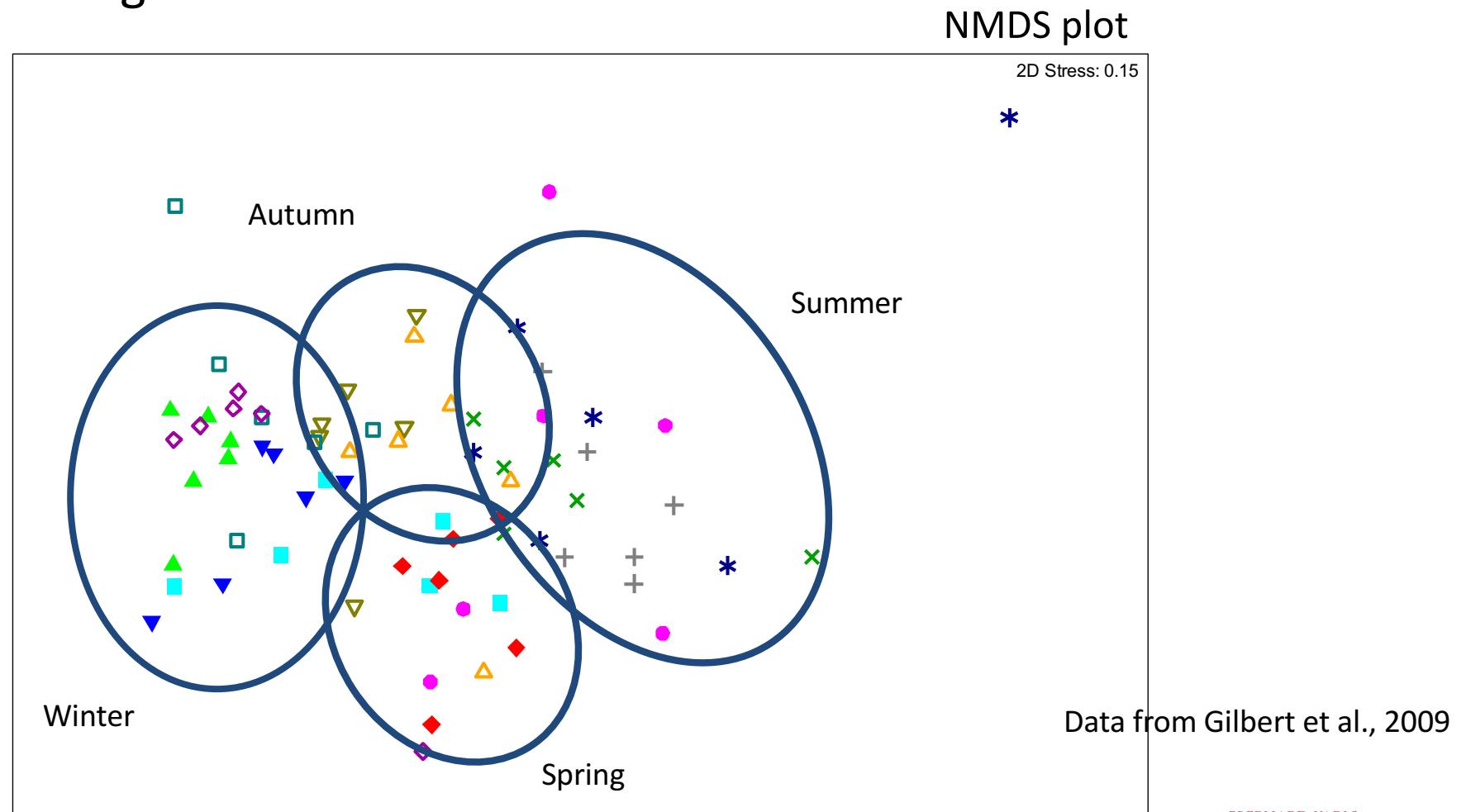
Comparison of eight marine metagenomes using six ecological measures



Correlation Study

Correlation between samples and different seasons

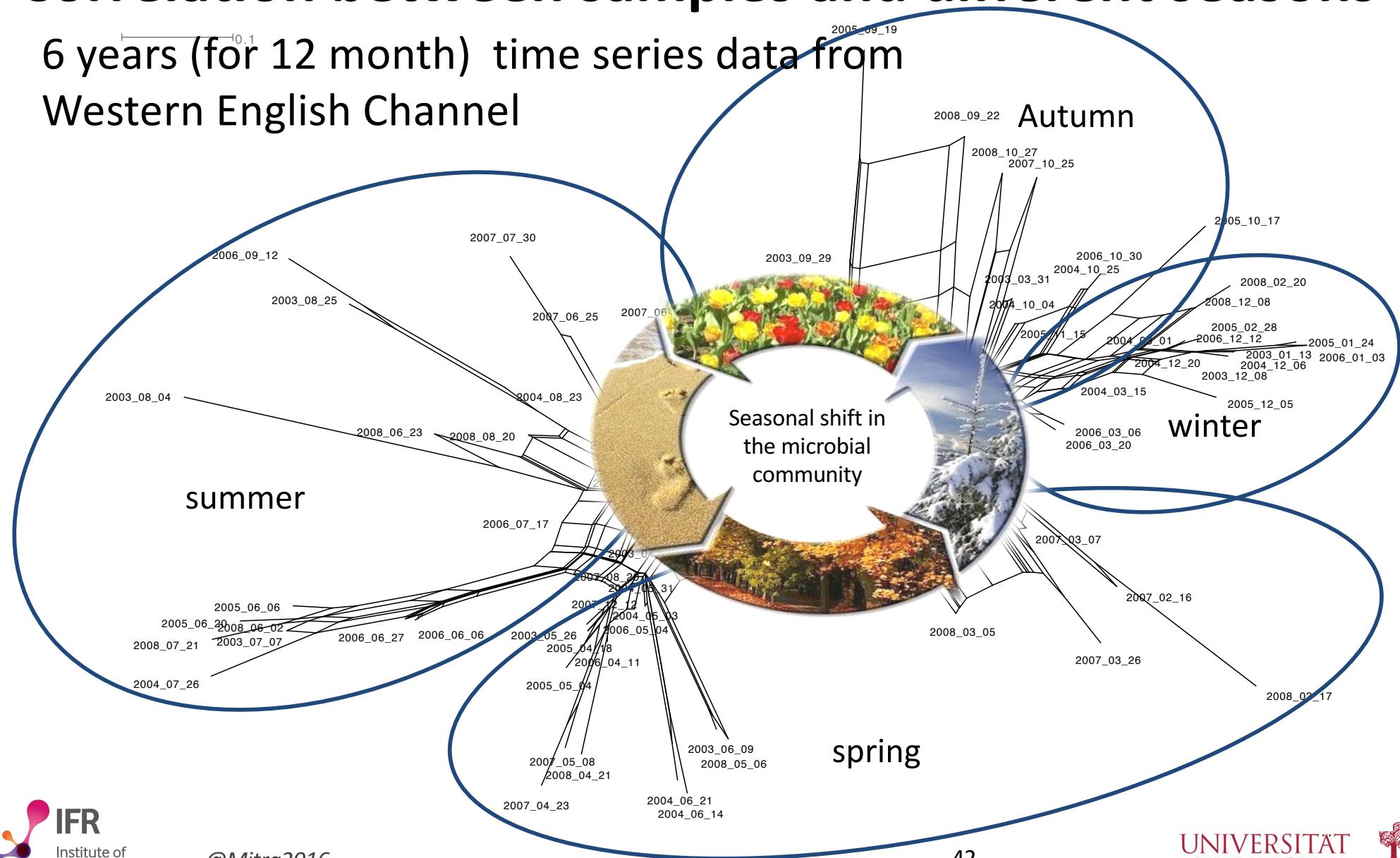
6 years (for 12 month) time series data from
Western English Channel



Correlation Study

Correlation between samples and different seasons

6 years (for 12 month) time series data from
Western English Channel

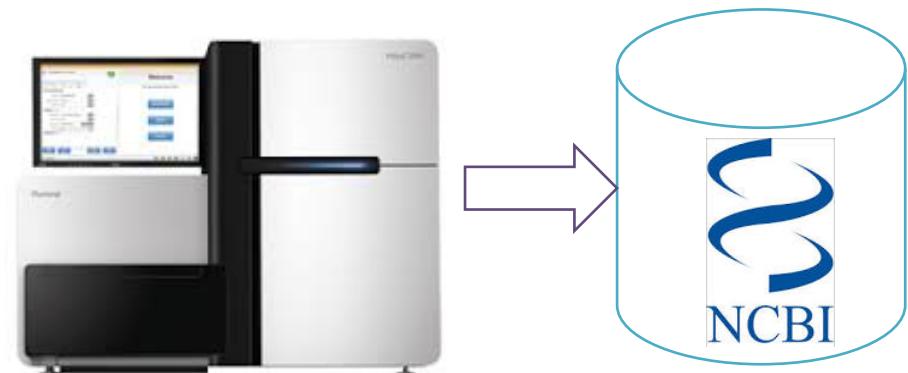


Bioinformatics and complexity



Bioinformatics

- Metagenomic shotgun data
 - One billion reads
 - 15 samples
- Alignment against NCBI-NR protein reference database
 - BLASTX on single server: ~15 years (*Altschul et al, 1990*)
 - RAPSearch2 on single server: ~3.5 months (*Zhao et al, 2012*)
- This takes too long...





DIAMOND

- ✓ DIAMOND is a “drop in” replacement for BLASTX
- ✓ Much, much faster...
- ✓ Very similar sensitivity to BLASTX
- ✓ Software available for download here:
<http://ab.inf.uni-tuebingen.de/software/diamond>

✓

Fast and sensitive protein alignment using DIAMOND

Benjamin Buchfink¹, Chao Xie^{2,3} &
Daniel H Huson^{1,2}

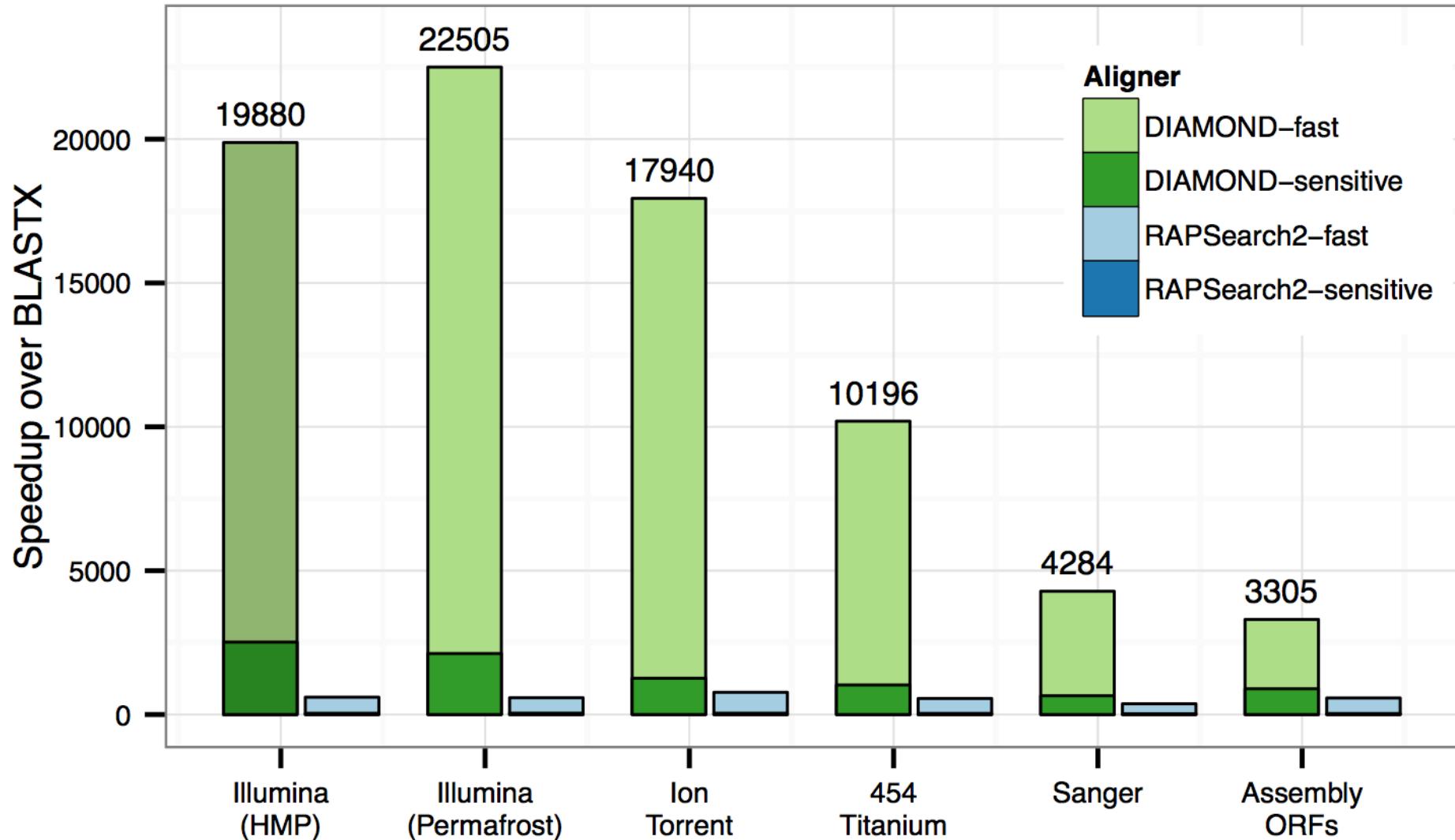
NATURE METHODS 2015



Performance

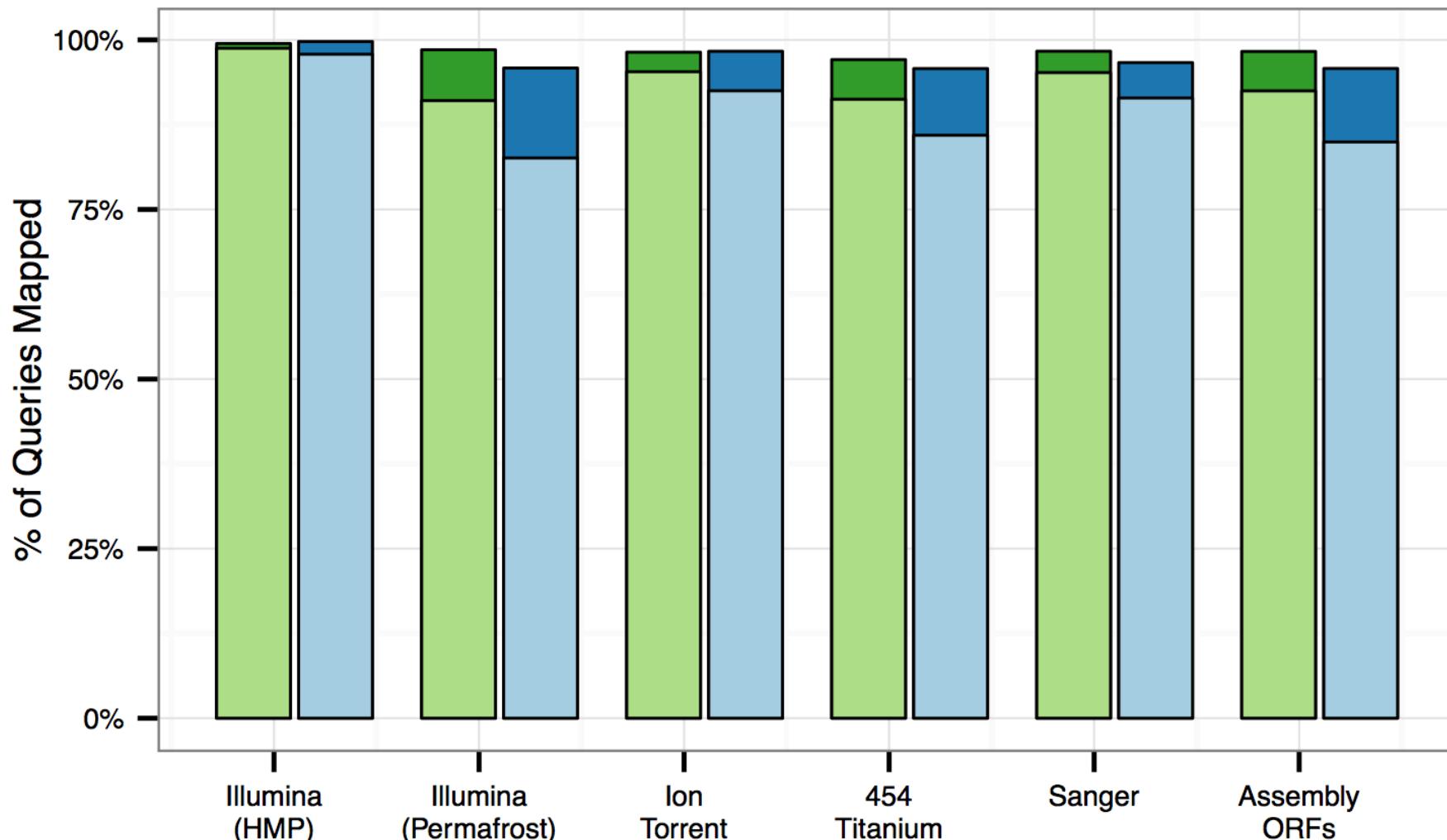
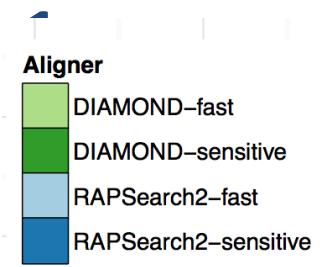
- DIAMOND is **20 000x** faster than BLASTX on Illumina HiSeq reads

DIAMOND - Speedup



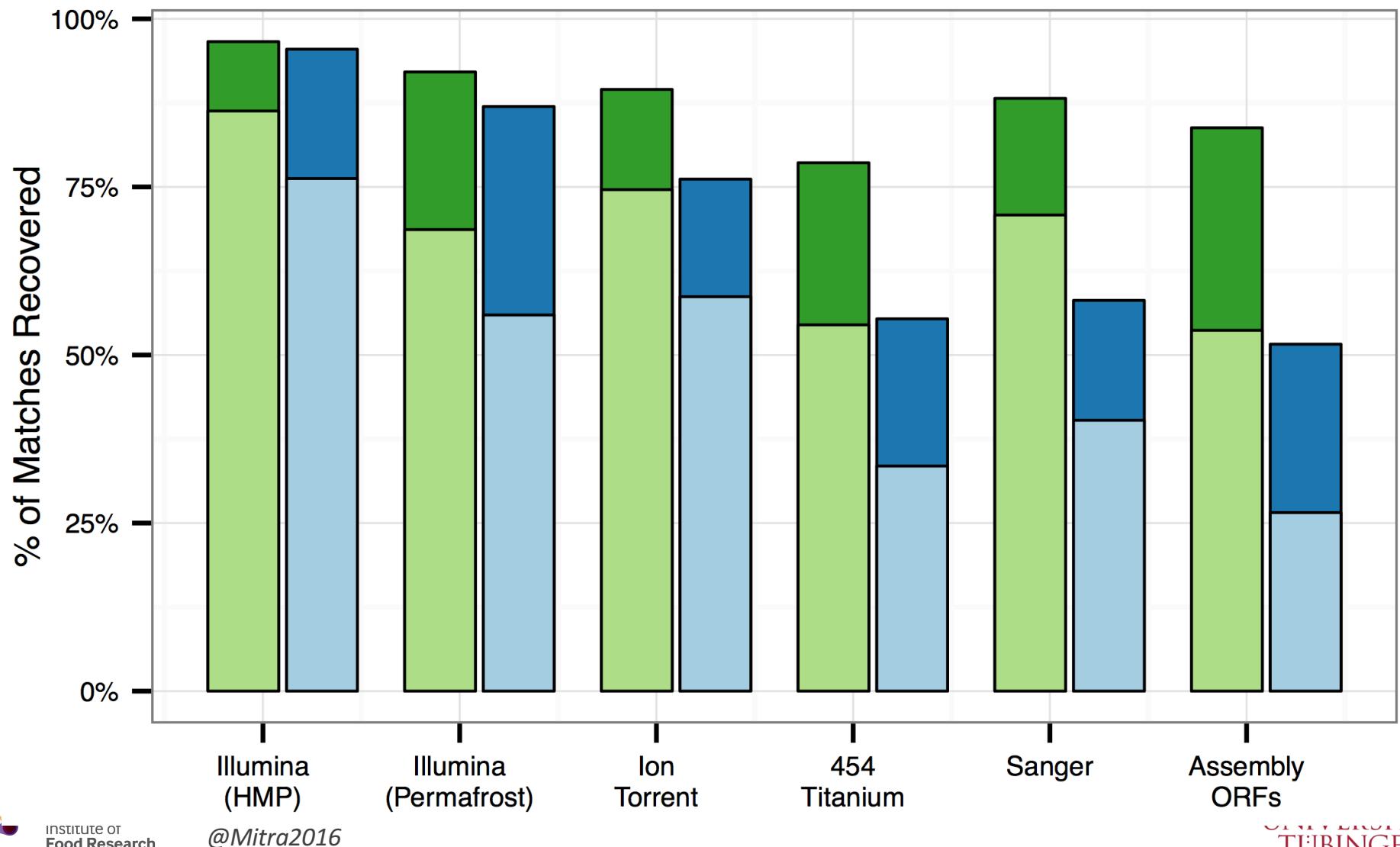
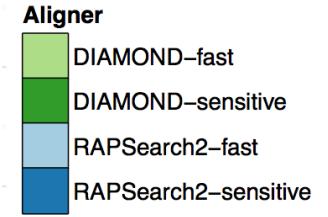
DIAMOND – Sensitivity

Percent of reads with a BLASTX match that also have a DIAMOND or RAPSearch2 match ($e\text{-value} \leq 10^{-3}$)



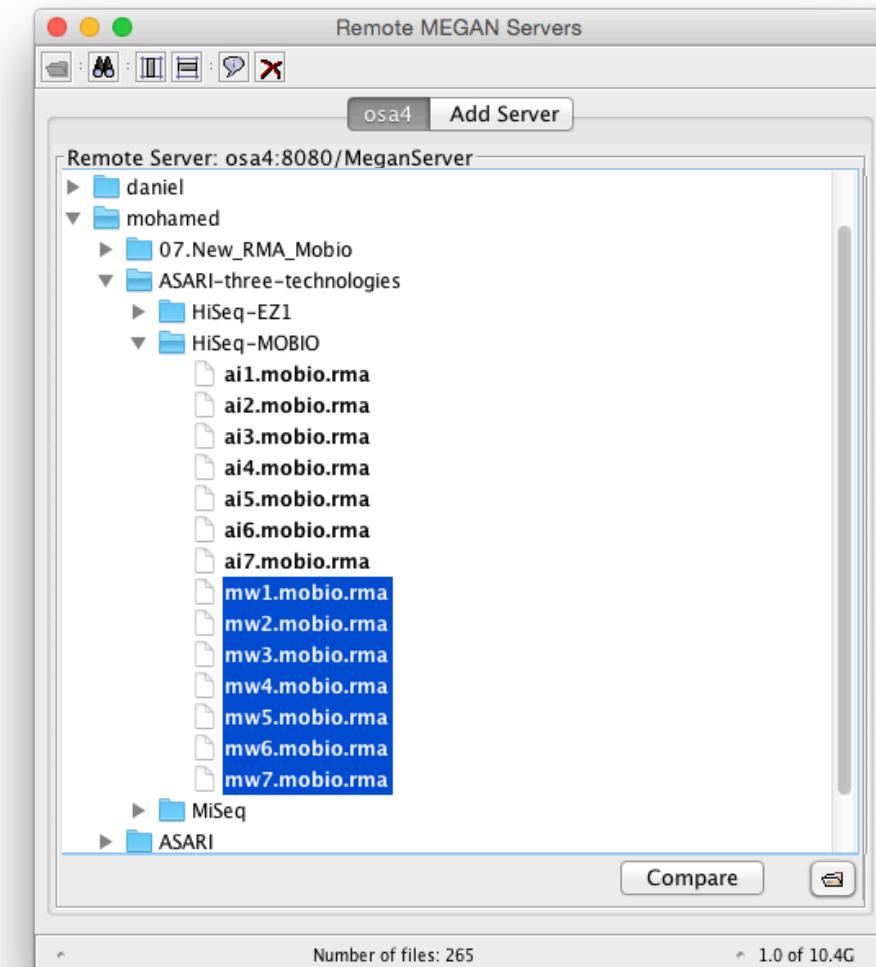
DIAMOND – Sensitivity 2

Percentage of BLASTX matches recovered
(≤250 matches per query, e-value ≤ 10^{-3})



Large Numbers of Large Samples

- Typical project has 100s of samples, each 10s of GBs
 - How to work with such data?
 - MeganServer software allows remote access from within MEGAN6



(Hans Ruscheweyh)

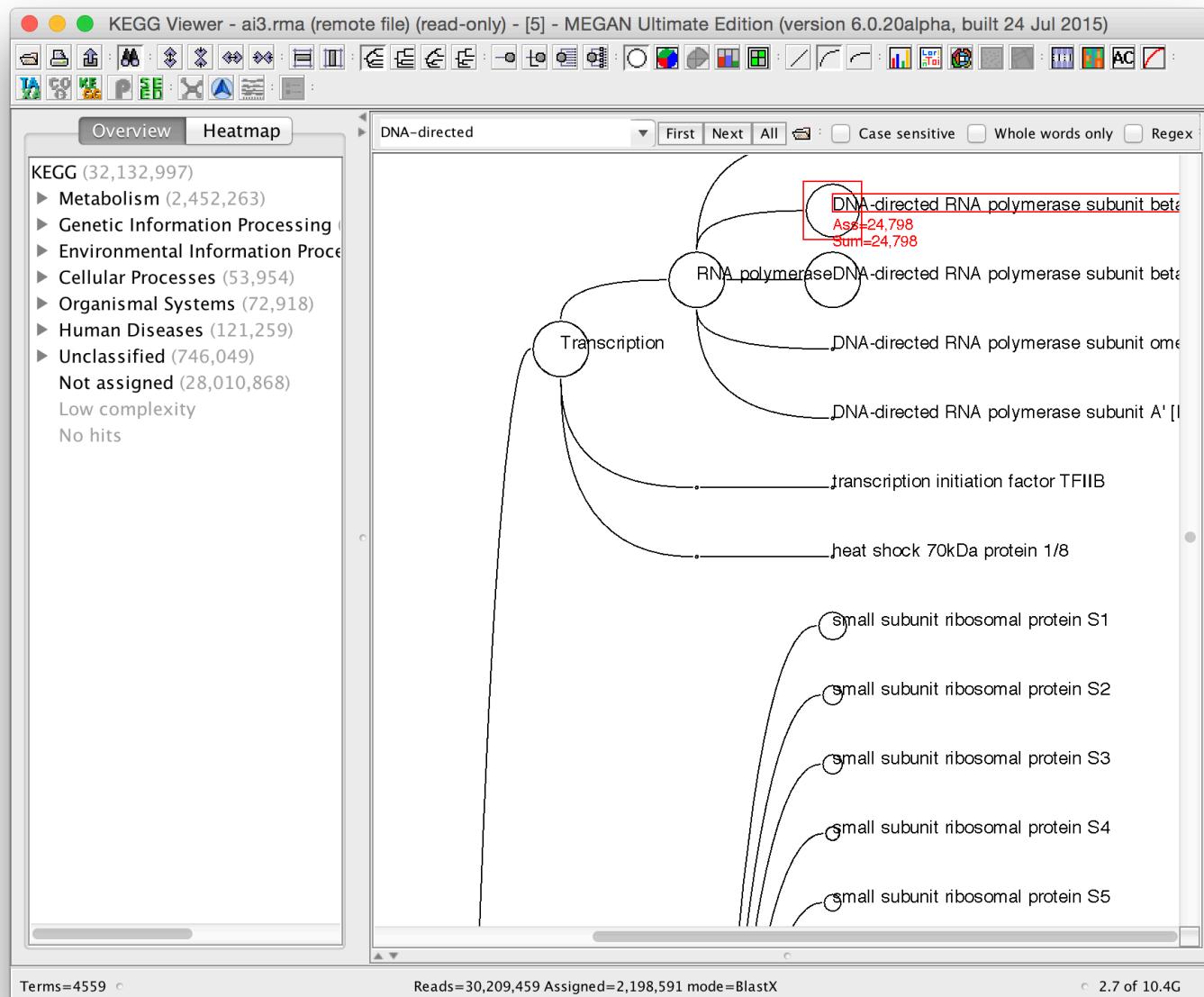
Metagenome Assembly

- Assembly of multiple genomes from metagenome data is very difficult
- Data reduction: first assemble, then BLAST
- Now: first (DIAMOND-)BLAST, then assemble

Gene-Centric Assembly

- Assemble gene-by-gene...
- Use DIAMOND to align reads to protein sequences
- Assemble all reads that align to a given gene
- Do this on-the-fly in MEGAN

Gene-Centric Assembly, rpoB



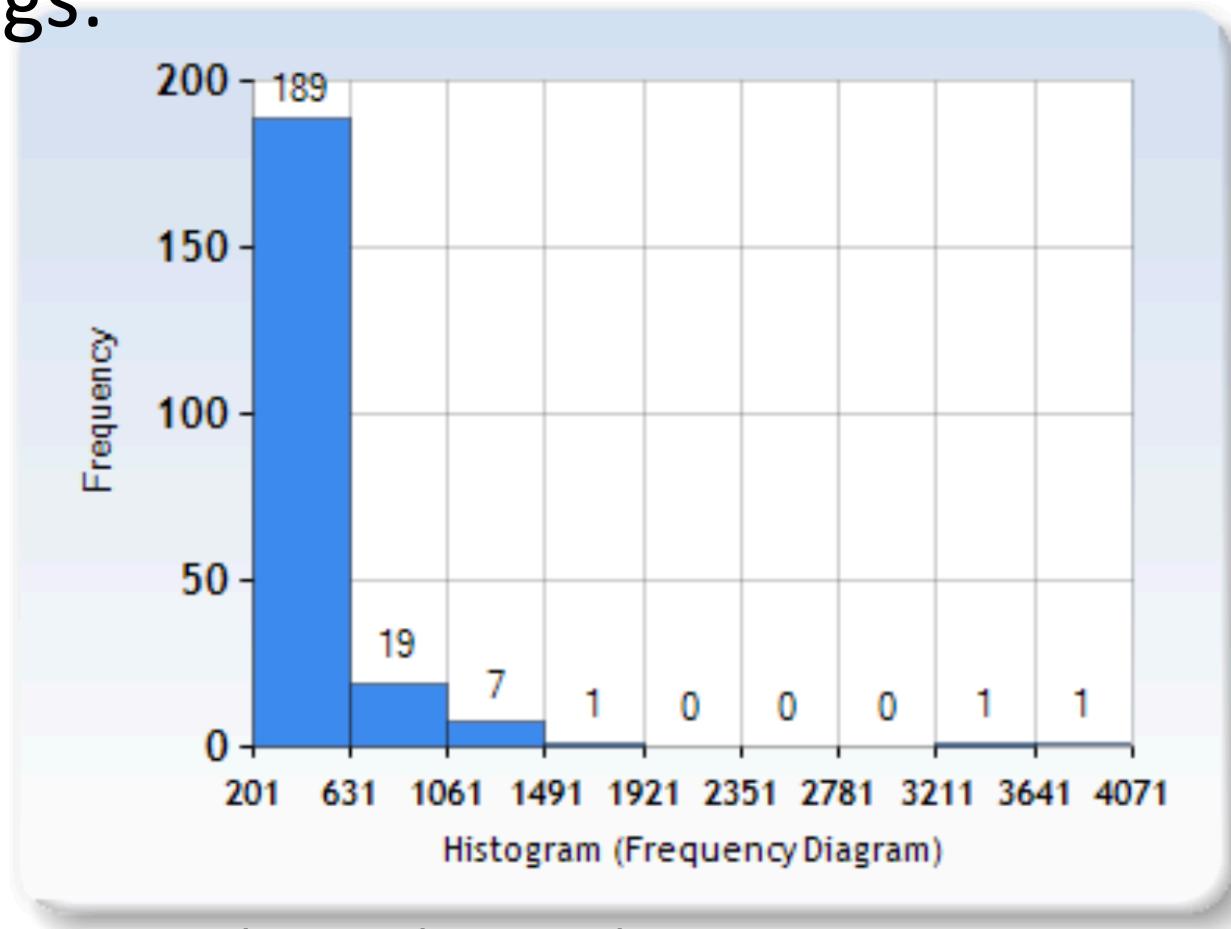
Human gut sample, 30 mio reads, 24,798 assigned to rpoB

Gene-Centric Assembly, rpoB



Gene-Centric Assembly, rpoB

- Gene-centric assembly of 24,798 reads, contigs:



+ 82 shorter than 200 bp

Will allow you to:

- Organize
- Visualize
- Interact
- Summarize
- Capture
- Compare

your metagenome data

<http://ab.inf.uni-tuebingen.de/software/megan6/>

Acknowledgements

- Tübingen, Germany
Prof. Daniel Huson
Hans-Joachim Ruscheweyh
Max Schubach
- SCELSE, Singapore
Prof. Stephan Schuster
Dr. Rohan Williams
- Karlsruhe, Germany
Dr. Bernhard Klar
- Chicago, USA
Dr. Jack Gilbert