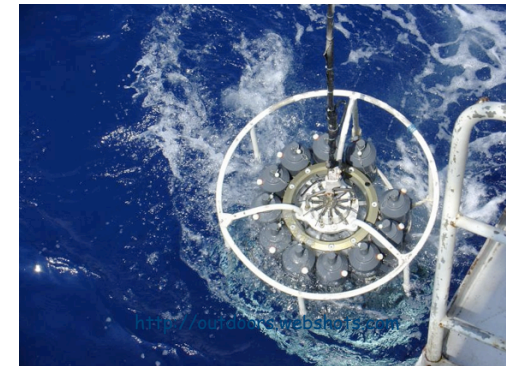
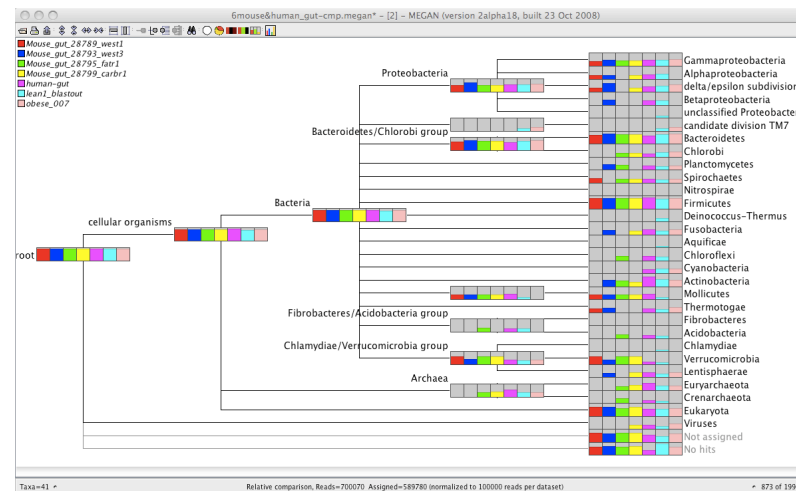




Computational Analysis of Metagenomes



Daniel H. Huson

PRIB Nijmegen 2010





Contents

- Genomics
- Sequencing
- Metagenomics
- Computational questions
- Outlook



Contents

- **Genomics**

- Sequencing

- Metagenomics

- Computational questions

- Outlook



Discovery of DNA

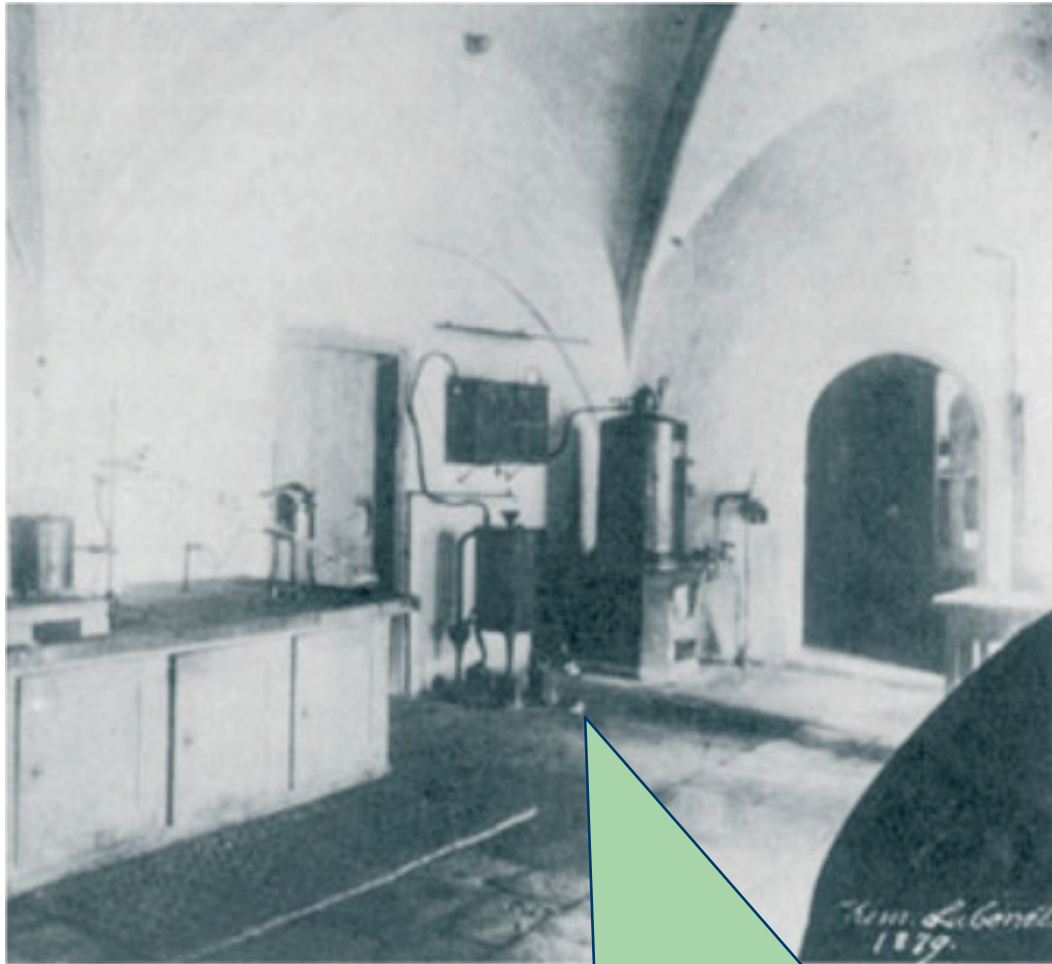
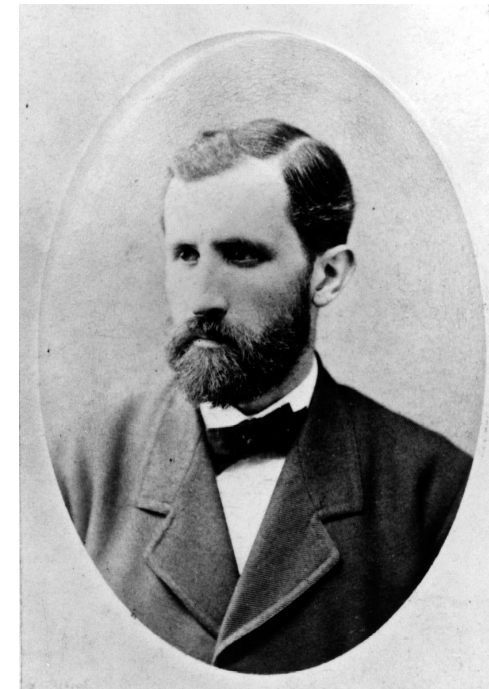
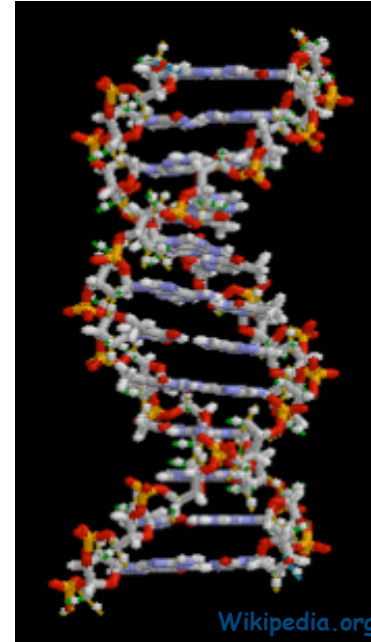


Photo courtesy of the University of Tübingen Library, Tübingen, Germany

**1869: Miescher discovered DNA
in the kitchen of Tübingen Castle**



**Friedrich Miescher
(1844-1895)**



...
A - T
C - G
G - C
T - A
A - T
A - T
...

1953 Watson and Crick

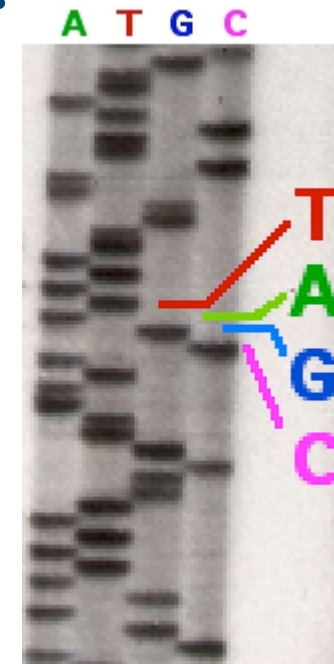
- The structure of DNA is a double helix
- It is *the order of the bases* along the molecule that contains heredity information

Sanger DNA Sequencing



1975 Frederick Sanger develops the “chain termination method” method for DNA sequencing

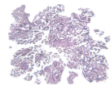
- Sanger sequencing basis of Genomics until 2005





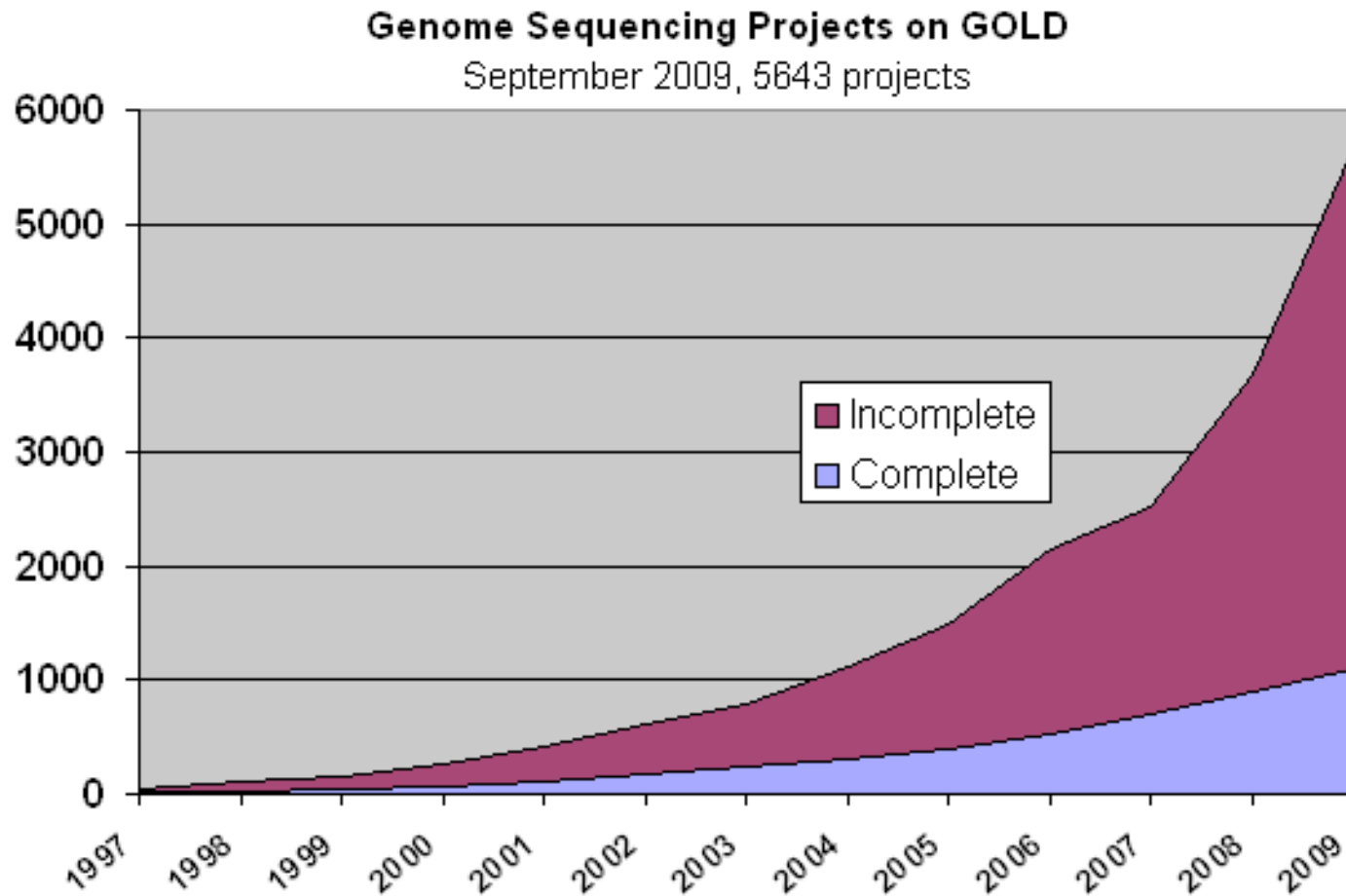
Genomics

- **Genomics** is the study of the genome sequence of individual organisms



- **Genome sizes:**
 - Bacteria: 1-10 million bases (Mb)
 - Drosophila: 140Mb
 - Human: 3 billion bases (Gb)

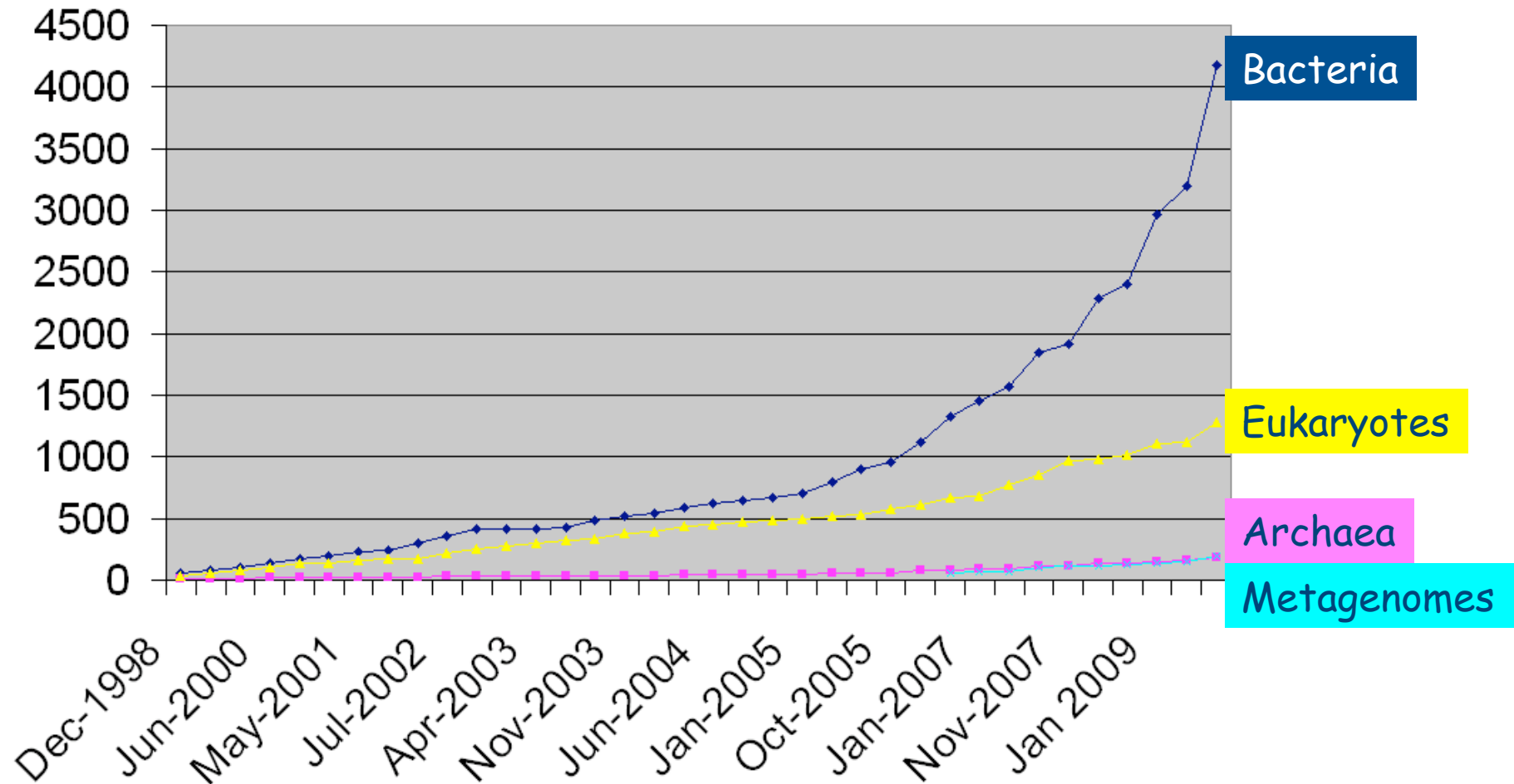
Sequencing of Genomes



GOLD: Genomes online database
www.genomesonline.org

Genome Projects by Groups

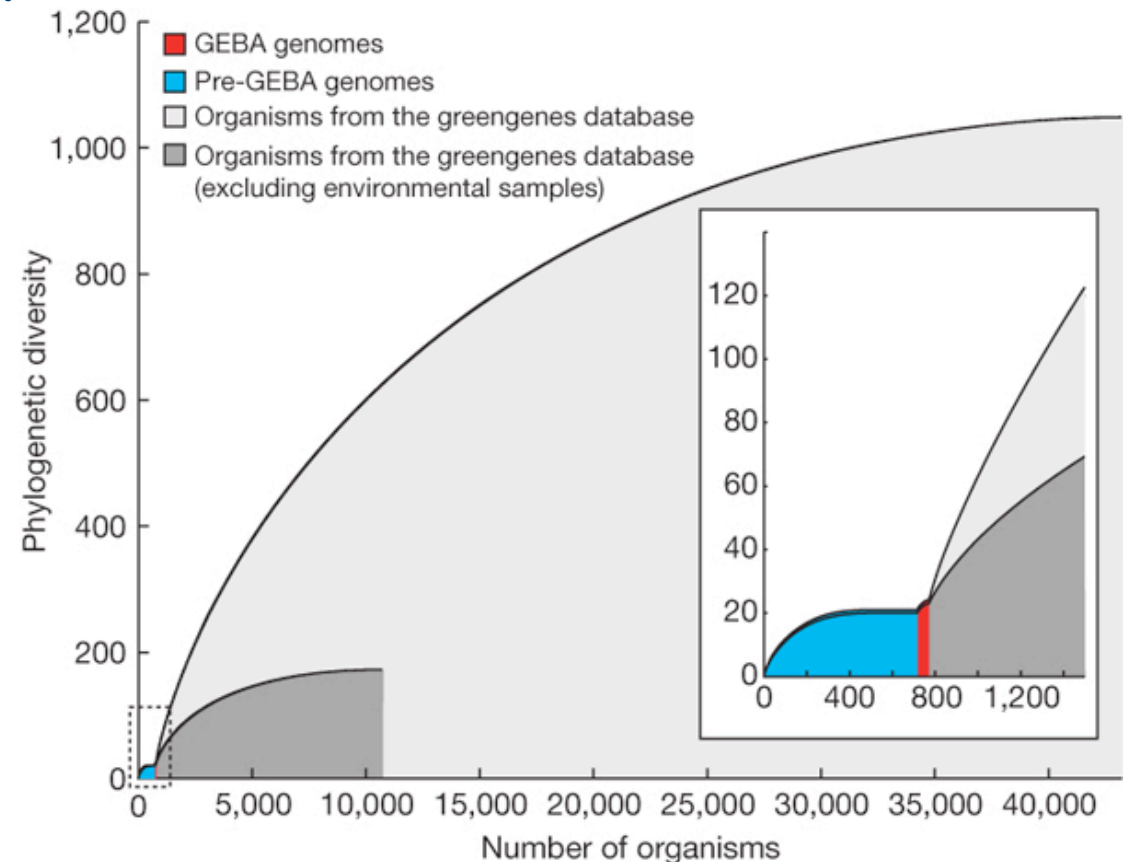
Genome Projects on GOLD according to Phylogenetic Groups ©
September 2009 - 5831 Projects



The GEBA Project

- A Genomic Encyclopedia for Bacteria and Archaea
 - JGI/DSMZ project

- Systematically
sequence
microbes from
underrepresented
clades



Dongying Wu *et al*, Nature, 2009

<http://www.jgi.doe.gov/programs/GEBA/>

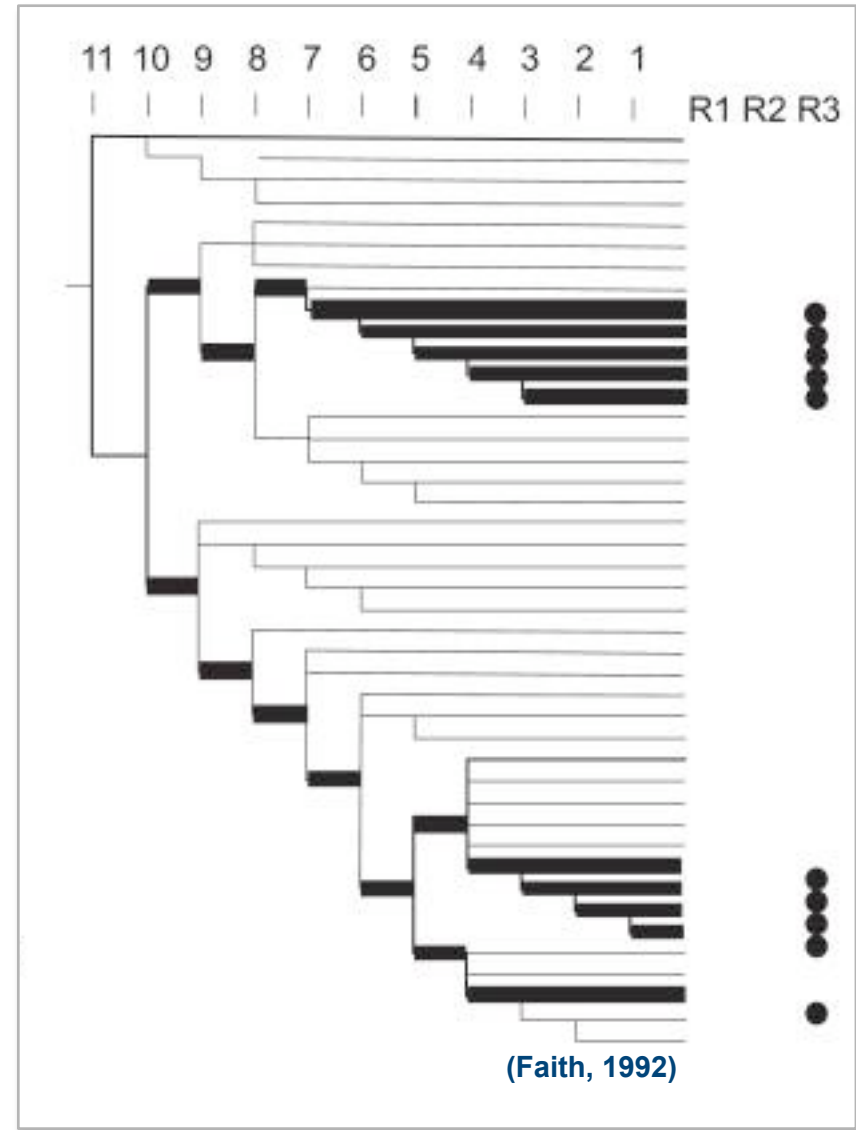


Phylogenetic Diversity (PD)

PD of a set of taxa:

- The sum of all branches on the phylogenetic tree that spans the set

(Faith, 1992)



From One to Many...

SCIENCE ONLINE | SCIENCE MAGAZINE HOME | SCIENCE NOW | NEXT WAVE | SCIENCE'S STKE | SCIENCE CAREERS

Science magazine | HELP | SUBSCRIPTIONS | FEEDBACK | SIGN IN

SEARCH | BROWSE | ORDER

The Sequence of the Human Genome

J. Craig Venter,^{1*} Mark D. Adams,¹ Eugene W. Myers,¹
Peter W. Li,¹ Richard J. Mural,¹ Granger G. Sutton,¹
Hamilton O. Smith,¹ Mark Yandell,¹ Cheryl A. Evans,¹
Robert A. Holt,¹ Jeannine D. Gocayne,¹ Peter Amanatides,¹
Richard M. Ballew,¹ Daniel H. Huson,¹ Jennifer Russo Wortman,¹
Qing Zhang,¹
Lin Chen,¹ M
Paul D. Thom
Catherine Nel
Joe Nadeau,⁵
Arnold J. Lev
Carolyn Slay
Arthur Delch



2001: THE Human Genome



Nature 409, 860 - 921 (2001) © Macmillan Publishers Ltd.

Initial sequencing and analysis of the human genome

International Human Genome Sequencing Consortium

The human genome holds an extraordinary trove of information about human development, physiology, medicine and evolution. Here we report the results of produce and make freely available a draft sequence of the human initial analysis of the data, describing some of the insights that ca

nature 15 February 2001



2008: 1000 Genomes Project...

Letter

Nature 463, 943-947 (18 February 2010) | doi:10.1038/nature08795; Received 11 August 2009; Accepted 6 January 2010

Complete Khoisan and Bantu genomes from southern Africa

See associated Correspondence: [Schlebusch, Nature 464, 487 \(March 2010\)](#), [Nature 464, 487 \(March 2010\)](#)

Stephan C. Schuster^{1,17}, Webb Miller^{1,17}, Aakrosh Ratan¹, Lynn P. Tomsho¹, Belinda Giardine¹, Lindsay R. Kasson¹, Robert S. Harris¹, Desiree C. Petersen², Fangqing Zhao¹, Ji Qi¹, Can Alkan³, Jeffrey M. Kidd³, Yazhou Sun⁴, Daniela I. Drautz¹, Pascal Bouffard⁴, Donna M. Muzny⁵, Jeffrey G. Reid⁵, Lynne V. Nazareth⁵, Qingyu Wang¹, Richard Burhans¹, Cathy Riemer¹, Nicola E. Wittekindt¹, Priya Moorjani⁶, Elizabeth A. Tindall^{2,7}, Charles G. Danko⁸, Wee Siang Teo^{2,7}, Anne M. Buboltz¹, Zhenhai Zhang¹, Qianyi Ma¹, Arno Oosthuisen⁹, Abraham W. Steenkamp¹⁰, Hermann Oostuisen¹¹, Philippus Venter¹², John Gajewski¹, Yu Zhang¹, B. Franklin Pugh¹, Kateryna D. Makova¹, Anton Nekrutenko¹,



Contents

- Genomics
- **Sequencing**
- Metagenomics
- Computational questions
- Outlook

Next-Generation Sequencing Technologies



- Fuelling a rapid growth of the number and size of sequencing projects



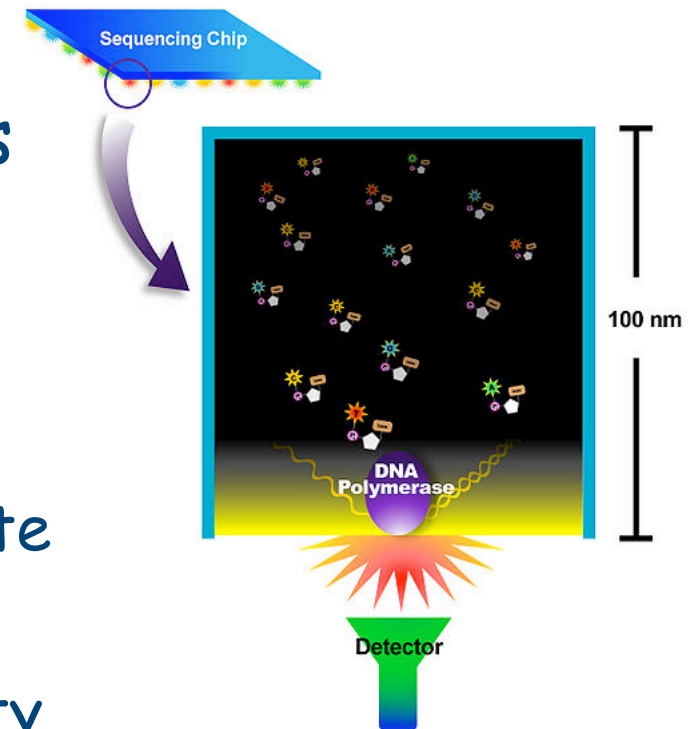
Advances in Sequencing Technologies

- First generation (Sanger sequencing):
 - 100kb/run, read length 1000bp, 500\$/Mb
- Second generation:
 - Roche/454: 450Mb/run, 400bp, 20\$/Mb
 - Illumina: 35Gb/run, 100bp, 0.50\$/Mb
 - SOLiD: 50Gb/run, 50bp, 0.50\$/Mb
 - Heliscope: 37Gb/run, 32bp, <0.50\$/Mb
- Third generation:
 - PacBio SMRT: 25Gb/run, >1000bp, ?\$/Mb
- Other:
 - Ion Torrent: uses ion sensor, <100,000\$

SMRT™ Sequencing

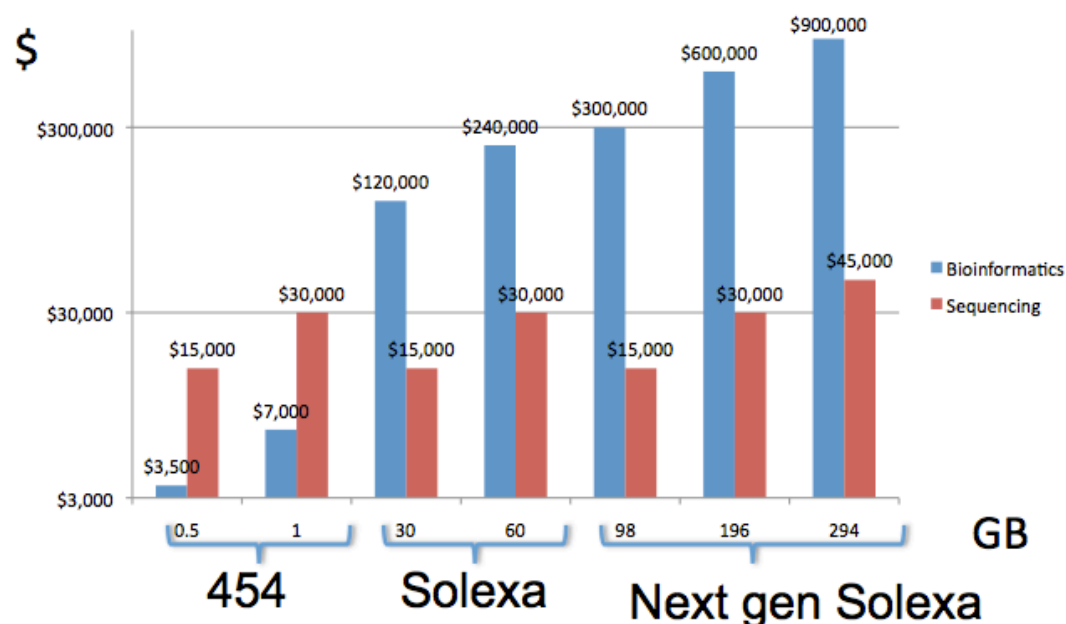
Single Molecule Real Time Sequencing

- Observes detached fluorescent dye molecules
- Three protocols:
 - **Linear sequencing:**
1kb reads, 10% deletion rate
 - **Circular sequencing:**
e.g. 200bp reads, high quality
 - **Strobe sequencing:**
e.g. 10 sections, each 100bp , each 500bp apart



Sequencing No Longer the Bottleneck...

If we simply ran BLASTX on EC2...



- 95GB == 195,600 node hours (on Nehalem 8core, 16GB),
- Illumina HiSeq2000 = 2x100GB/run
- cost is purely BLAST, no storage or transfer cost
- values are in Amazon EC2 (from Wilkenina et al, *IEEE Cluster09*)
- note: 10x or 100x improvements over BLASTX will help, but not solve
- prices from mid 2009



This slide kindly provided by Folker Meyer (Argonne National Labs)



Contents

- Genomics
- Sequencing
- **Metagenomics**
- Computational questions
- Outlook



How Many Species?



Major unsolved question:

- Number of species on Earth?
- Cannot be answered even to within several orders of magnitude
- Some estimations
 - 3-50 million species of arthropods
 - 1-100 million species of nematodes



www.ucmp.berkeley.edu/arthropoda/arthropoda.html

Once the diversity of the microbial world is catalogued, it will make astronomy look like a pitiful science

- Julian Davies, Professor Emeritus, Microbiology and Immunology, UBC



Identified Modern Species

~1.7 million named species

- 287,655 plants, including:
 - 15,000 mosses
 - 13,025 ferns
 - 980 gymnosperms
 - 199,350 dicotyledons
 - 59,300 monocotyledons
- 74,000-120,000 fungi
- 10,000 lichens
- 5,700 prokaryotes
- ~1,250,000 animals, including:
 - 1,190,200 invertebrates:
 - 950,000 insects
 - 70,000 mollusks
 - 40,000 crustaceans
 - 130,200 others
 - 58,808 vertebrates:
 - 29,300 fish
 - 5,743 amphibians
 - 8,240 reptiles
 - 10,234 birds
 - 5,416 mammals

Source: <http://en.wikipedia.org/wiki/Biodiversity>

Sequences for ~200,000



Metagenomics

- “The study of the DNA of uncultured organisms”
- > 99% of all microbes cannot be cultured
- A genome:
 - Entire genetic information of a single organism
- A metagenome:
 - Entire genetic information of a community of organisms



www.innovations-report.de

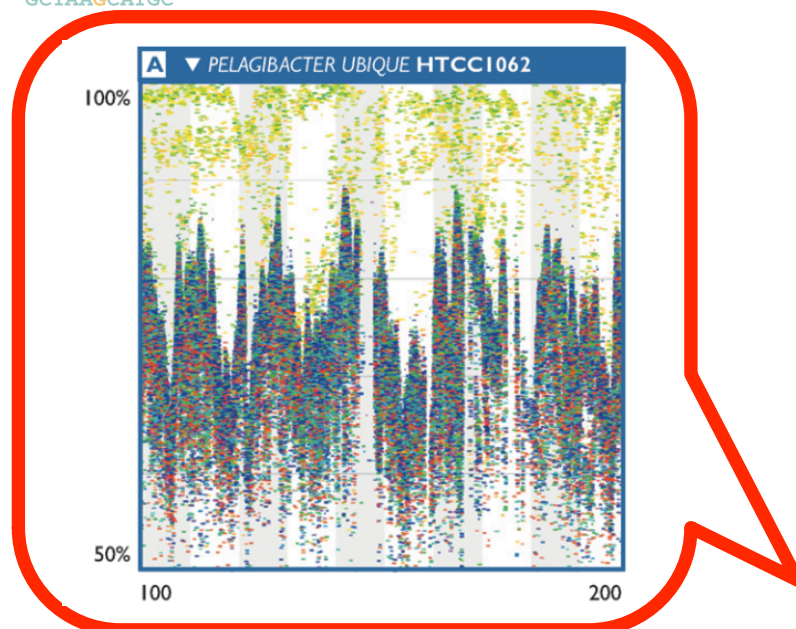


Typical Sources of Metagenomes

- Soil samples
- Sea water samples
- Seabed samples
- Air samples
- Medical samples
- Ancient bones
- Human microbiome

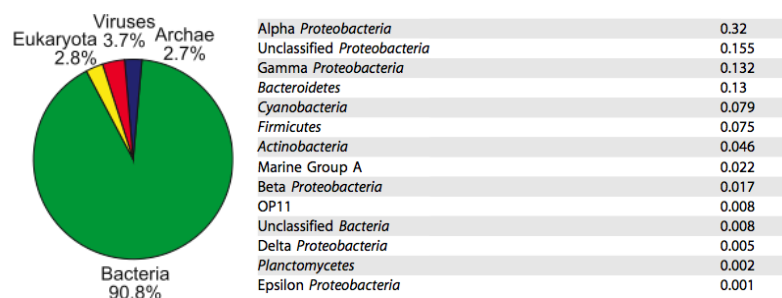


Global Ocean Sampling Expedition



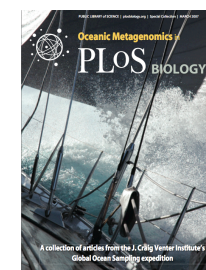
Rusch *et al.* (2007):

- 41 samples
- Size filtered 0.1-0.8 μ m
- Sanger sequencing
 - 7.7 million reads
 - length ~822bp
 - ~ 5.9Gb sequence
- Low abundance of *clonal* organisms



Yooseph *et al.* (2007):

- 6 million proteins
 - linear rate of discovery





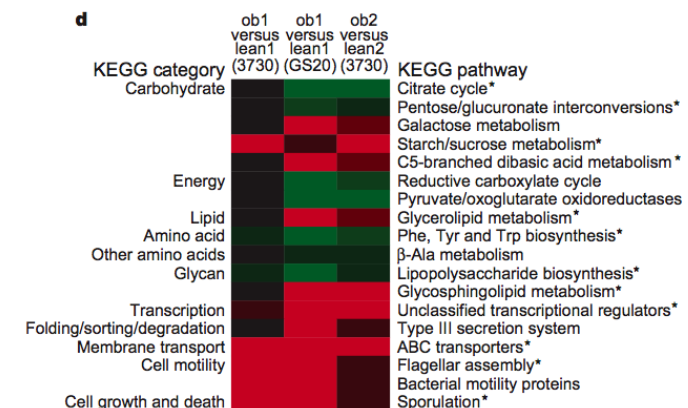
Gut Microbiota



- 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

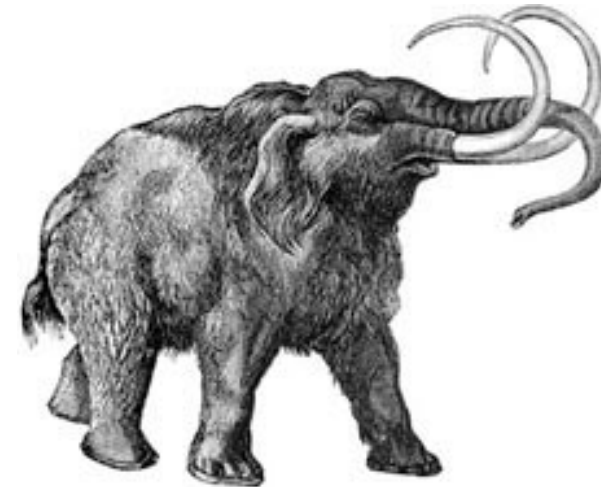
• Obesity-associated gut microbiome

- Change in relative abundance of Bacteroidetes and Firmicutes
- Change in functional capacity (toward energy harvesting)

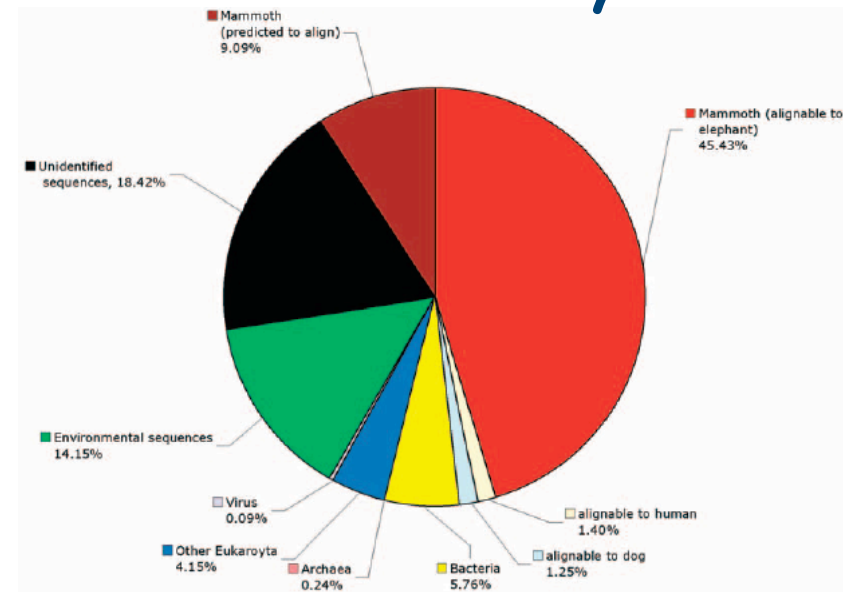


Mammoth Project

- DNA collected from permafrost mammoth (28,000 years old)
- DNA extracted from 1 gram of bone
- 454 sequencing:
 - ~302,000 reads
 - ~95 bp length
- > 50% mammoth



Taxonomic analysis



REPORTS

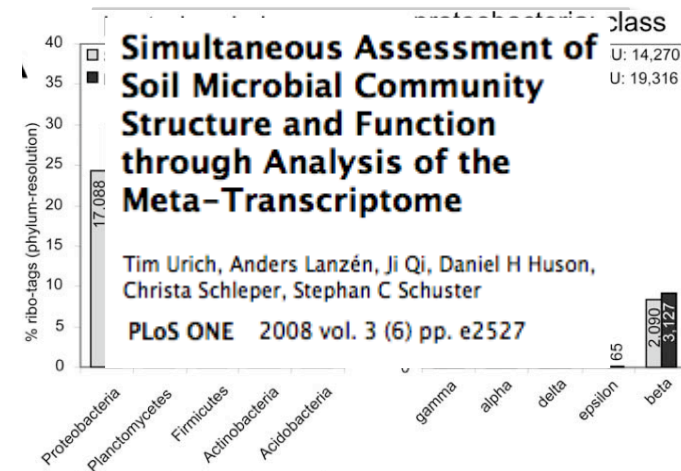
Metagenomics to Paleogenomics: Large-Scale Sequencing of Mammoth DNA

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

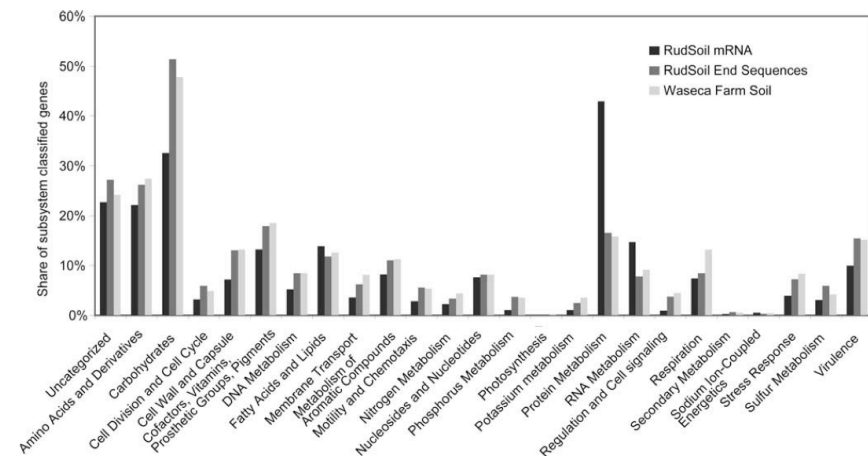


"Meta Transcriptomics" of Soil

- Urich *et al* (2008):
 - RNA randomly reverse transcribed into cDNA
 - No PCR or cloning
 - 454 sequencing:
 - ~ 250,000 sequences
 - ~ 98 bp length
 - RNA types:
 - ~ 75% rRNA tags
 - ~ 8% mRNA tags
 - ~ 17% unassigned



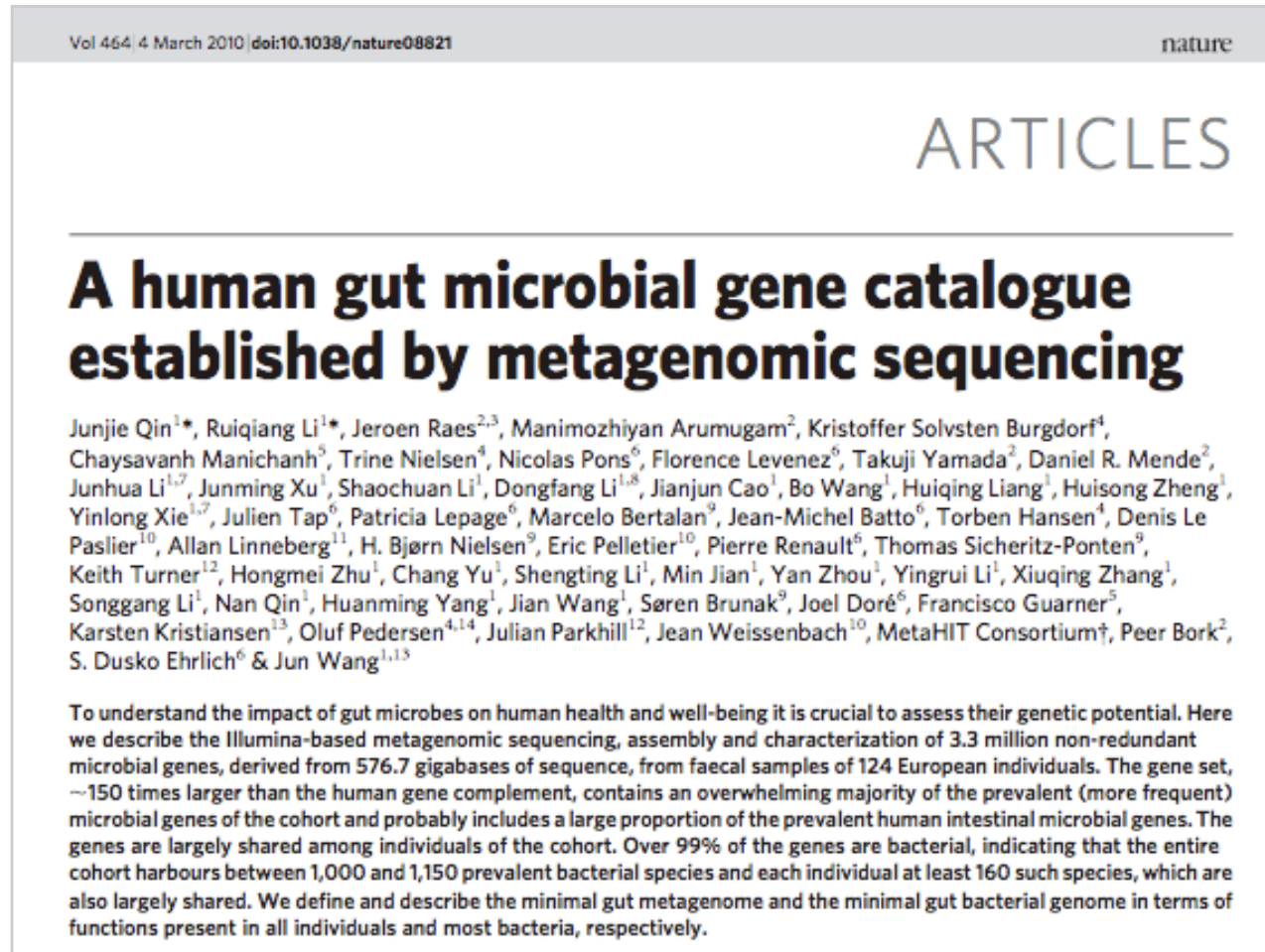
rRNA analysis



mRNA analysis



Large-Scale Human Gut Analysis



- 576Gb of sequence from 124 individuals



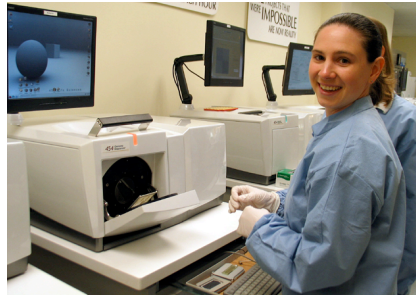
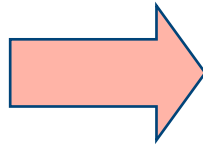
Contents

- Genomics
- Sequencing
- Metagenomics
- **Computational questions**
- Outlook

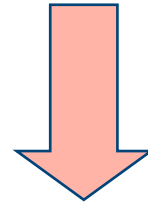
Metagenome Analysis



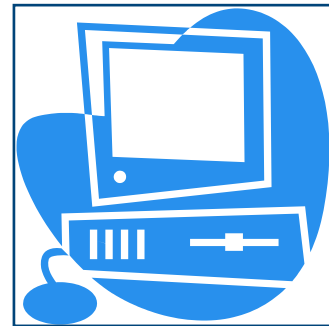
Environmental
sample



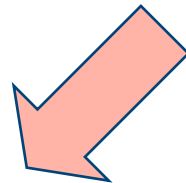
High-throughput
DNA sequencing



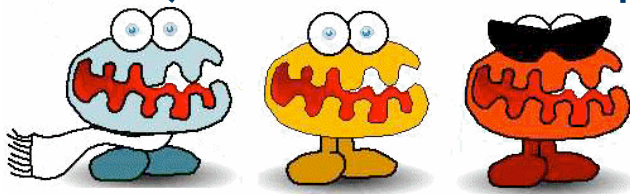
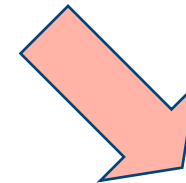
10 million sequences



Basic computational
analysis



10 000
hours



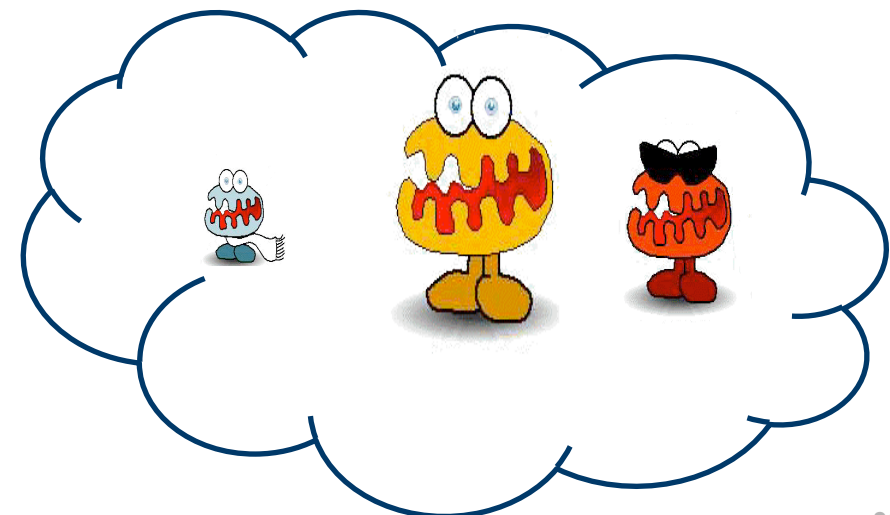
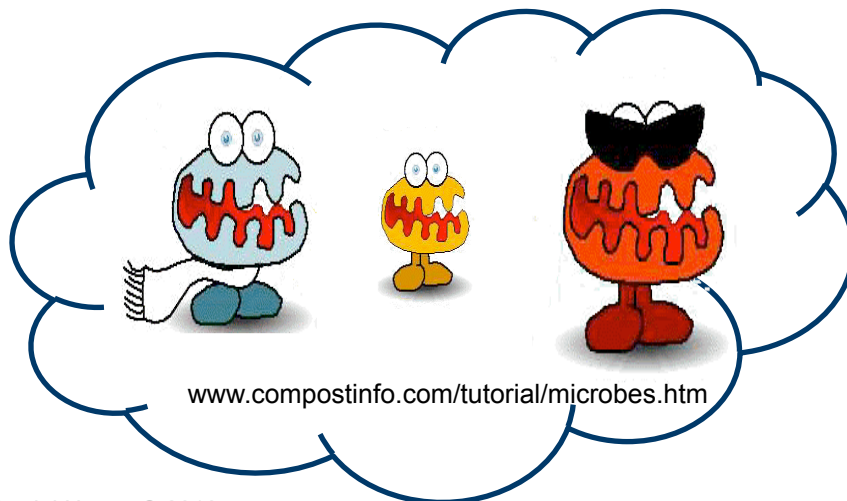
Q1: Who is out there?

Daniel Huson © 2010



Q2: What are they doing?

Q3: How Do They Compare?





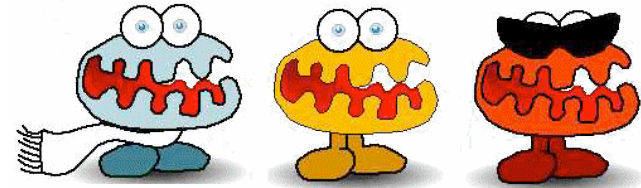
Additional Questions

- How to cluster reads by relatedness (using machine learning techniques)?
- How to assemble metagenome data?
- Gene prediction?
- Faster sequence comparison
- etc

Three Basic Computational Questions

- Who is out there?

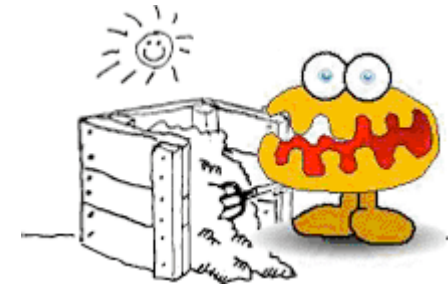
- Types of organisms
- In what proportions?



www.compostinfo.com/tutorial/microbes.htm

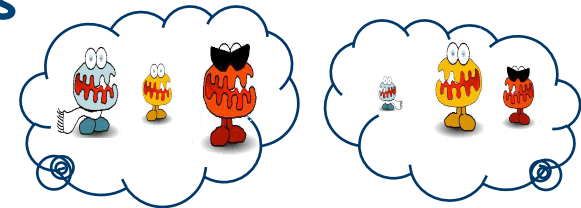
- What are they doing?

- 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

Three Basic Computational Questions

- Who is out there?

- Types of organisms
- In what proportions?



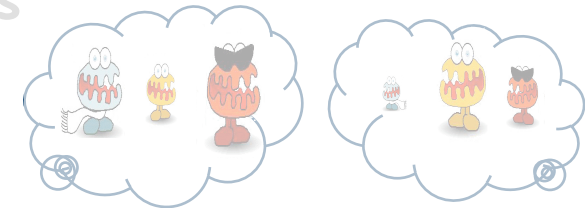
- What are they doing?

- 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?

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



Who is Out There?

Main tool: *Similarity search*

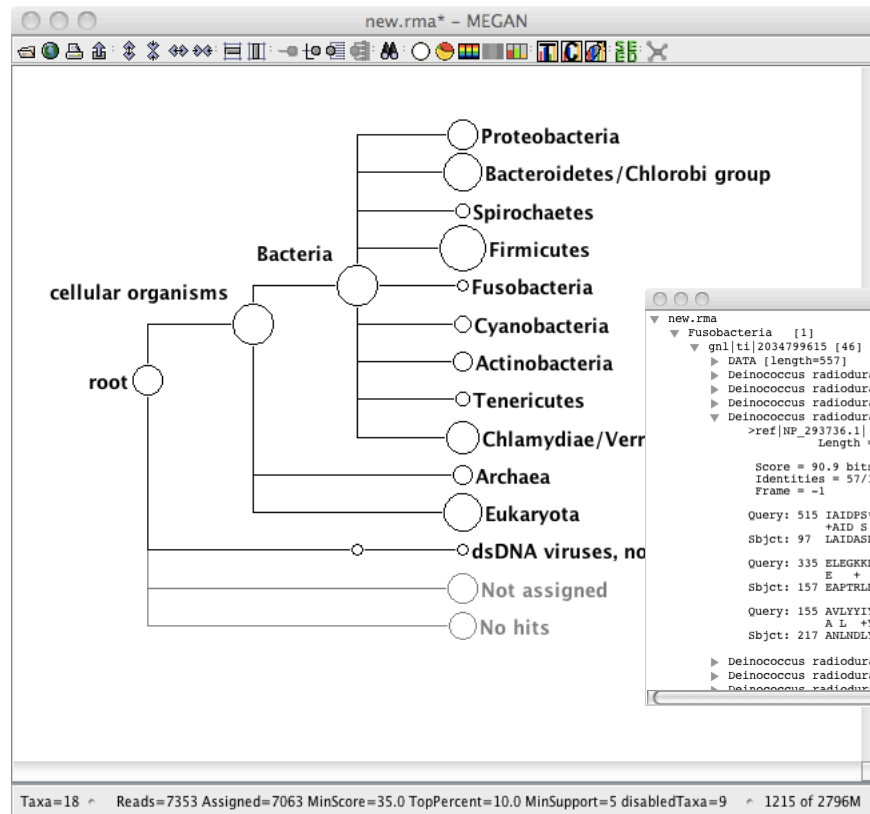
For every DNA (or cDNA) read:

- Find significant matches to sequences in a reference database
- Use matches to place read in NCBI taxonomy

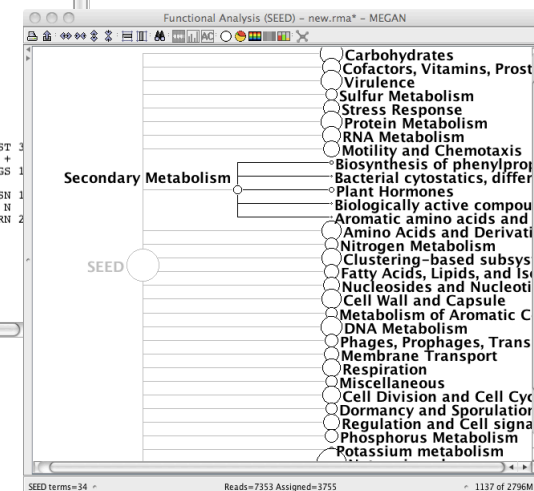
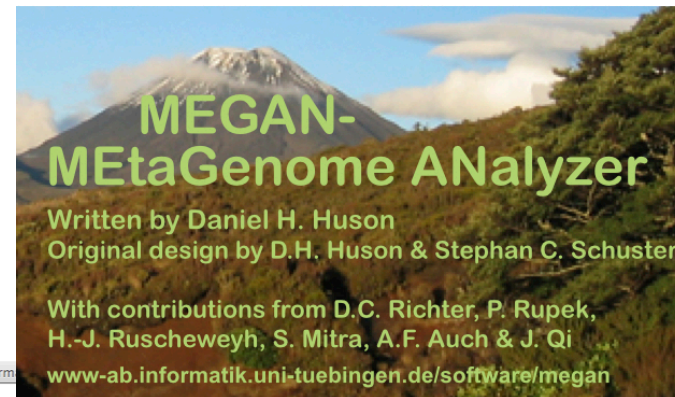
MEGAN - MEtaGenome ANalyzer



Huson et al, 2007

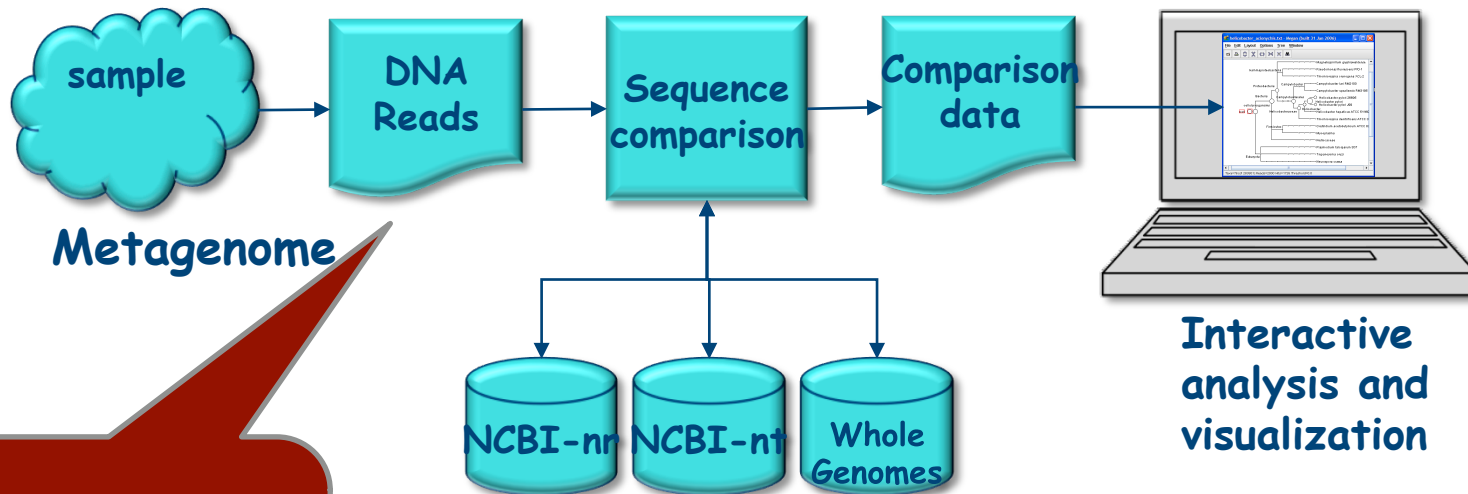


```
new.rma
  Fusobacteria [1]
    gnl|ti|2034799615 [46]
      DATA [length=557]
        Deinococcus radiodurans R1 score=90.9
        Deinococcus radiodurans R1 score=90.9
        Deinococcus radiodurans R1 score=90.9
        Deinococcus radiodurans R1 score=90.9
        >ref|NP_293736.1| putative transposase [Deinococcus radiodurans R1]
          Length = 400
          Score = 90.9 bits (224), Expect = 4e-17
          Identities = 57/149 (38%), Positives = 78/149 (52%)
          Frame = -1
          Query: 515 IAIDFS*VSKAGKTAHIGRFWSGCASAVKHGLEILGIAVIDADIRDAHMLRAVOTLNST
            +AID S KAG+ TAH+G FW+GCA+ + G+E A+ID R A+ + QTL +
            Sbjet: 97 LAIDASPHRAGHTAHLSFWNGCAARTERGIEQSCCALIDVQHRQALTVDRQTLTGS
          Query: 335 ELEGKFTLNQWYLSVLKTYRTOLLKITSLLVADAFAFVLPFVEGLKEIGFSLISRLRN
            E + VL RT + +VAD ++ F VE + G ISRL N
            Sbjet: 157 EAPTRLEQXADQLDDVLLDRTVQQLDLAAVADGNYAKEFIVETVTHGHLFISRLPN
          Query: 155 AVLYIIEGPRTKGRPRKTKDGKIDFSN 69
            A L +Y G +RGR K DGR+DFS+
            Sbjet: 217 ANLNDLYTGEHPRRGRKKKFDGKVDPSD 245
        Deinococcus radiodurans R1 score=90.9
        Deinococcus radiodurans R1 score=90.1
        Deinococcus radiodurans R1 score=88.4
```



- Interactive tool for metagenomic analysis
(Version 4, to be released Nov 2010)

Metagenomics Pipeline



Similar for

- metatranscriptomics
- metaproteomics
- amplicon sequencing

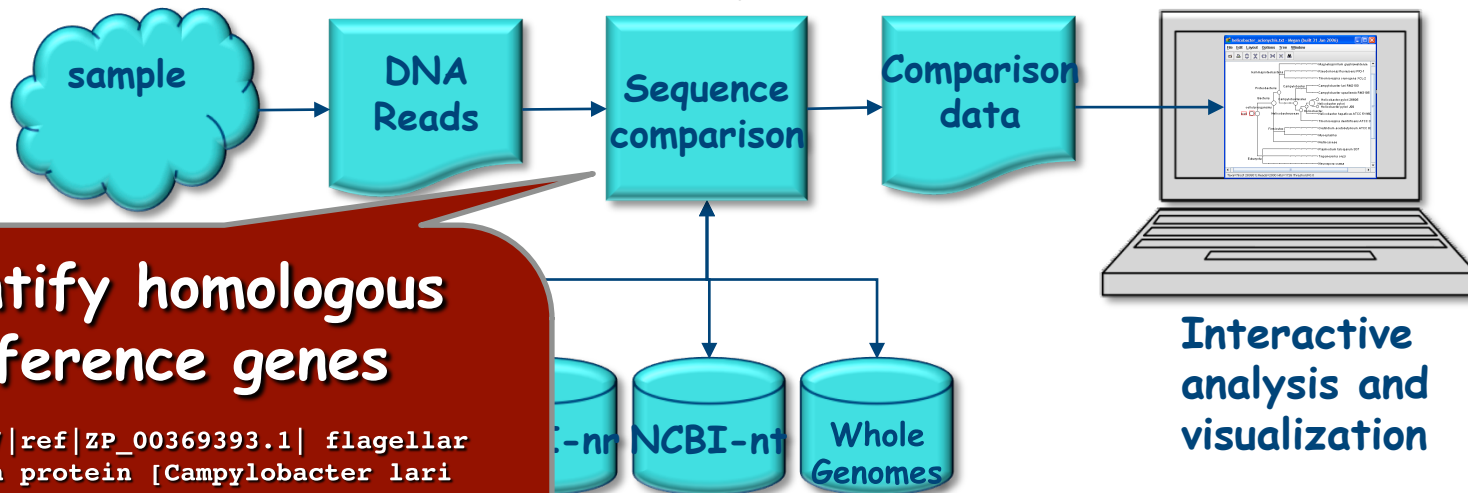
Reference databases

Identifying Taxa and Genes

Metagenome analysis

Basic idea: compare reads against references sequences of known species and/or function

**BLASTX
against
NCBI-NR**



Identify homologous reference genes

>gi|57241447|ref|ZP_00369393.1| flagellar motor switch protein [Campylobacter lari RM2100]

Score = 33.9 bits (76), Expect = 1.8
Identities = 13/26 (50%), Positives = 19/26 (73%)

Query: 79 LMFVFDDLATVEENGIREIINRADKK 2
LMF FDD++ + N IRE++ ADK+
Sbjct: 243 LMFTFDDISQLSTNAIREVLKAADKR 268

Reference databases

Sequence Comparison

- DNA Read



- Align to reference sequences, e.g.
BLASTX against NR database:

```
>gi|57241447|ref|ZP_00369393.1| flagellar motor switch protein
Campylobacter lari RM2100
Score = 33.9 bits (76), Expect =0.01
Identities = 13/26 (50%), Positives = 19/26 (73%)

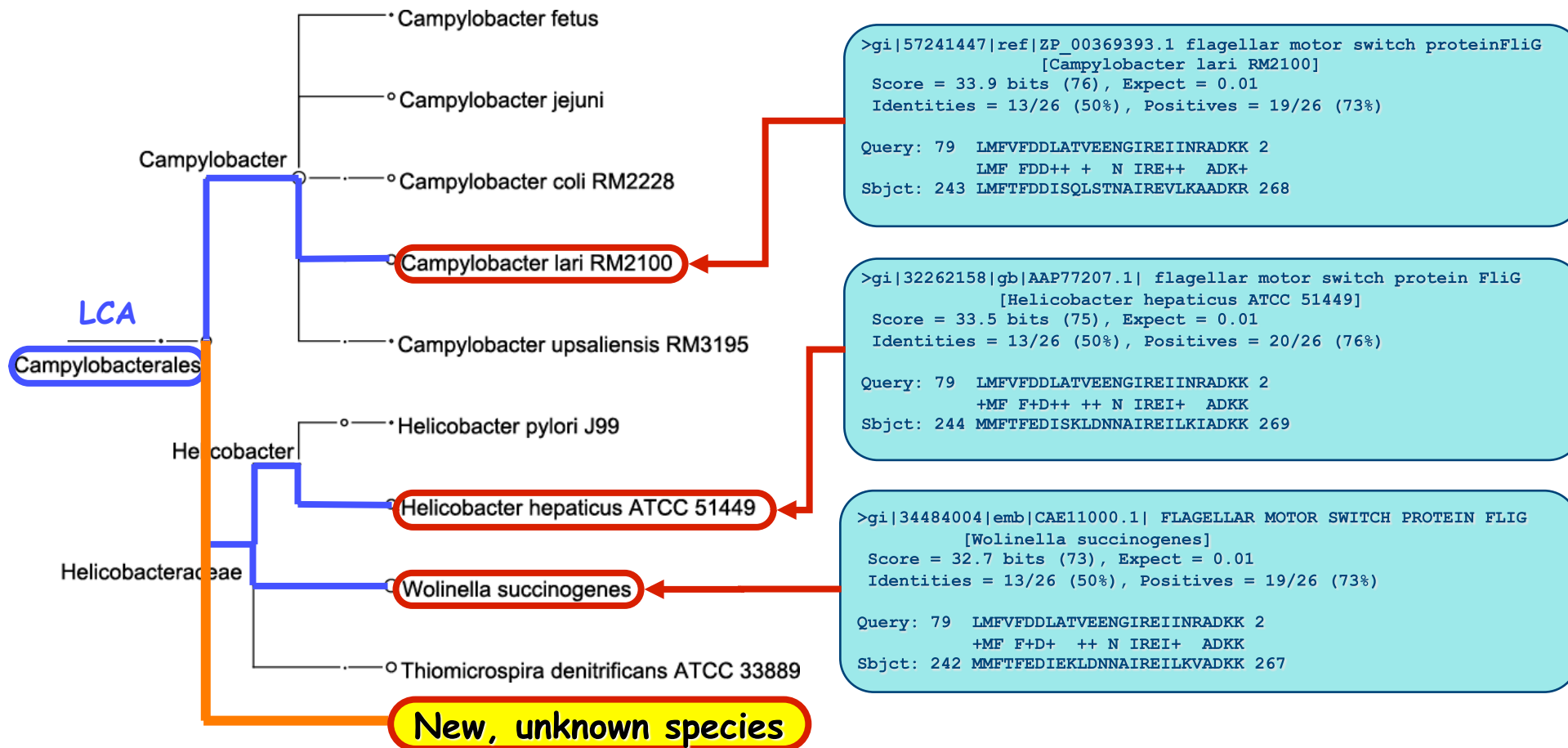
Query: 79  LMFVFDDLATVEENGINEIRELINRADKK 2
          LMF FDD++ + N IRE++ ADK+
Sbjct: 243 LMFTFDDISQLSTNAIREVLKAADKR 268
```

- Indicates gene content:

Campylobacter lari RM2100

Taxonomic Placement Using LCA

A read will often match more than one database entry:



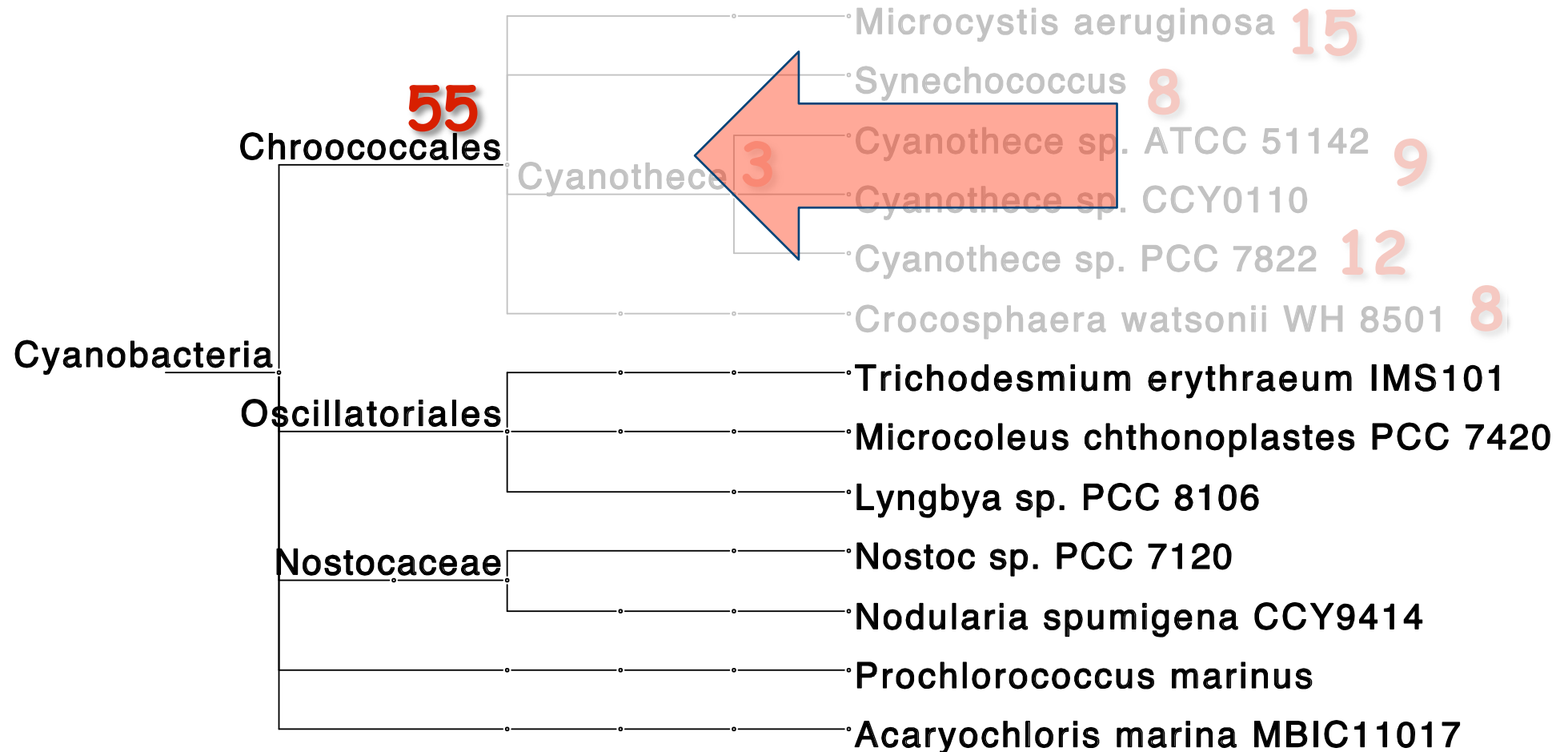
LCA approach: Assign read to LCA of hits in taxonomy



Taxonomic Placement Using LCA

- For each DNA read:
 - Determine which gene sequences it matches
 - Corresponding species are assumed to contain the gene
 - Place read on the LCA of species
- Is placement by **gene content** or phylogenetic footprint
- Robust against false positive placements
- Robust against (known cases) of horizontal gene transfer

Minimum Support Filter



- Require at least e.g. 50 reads on a node

Taxonomic Analysis using NCBI Taxonomy

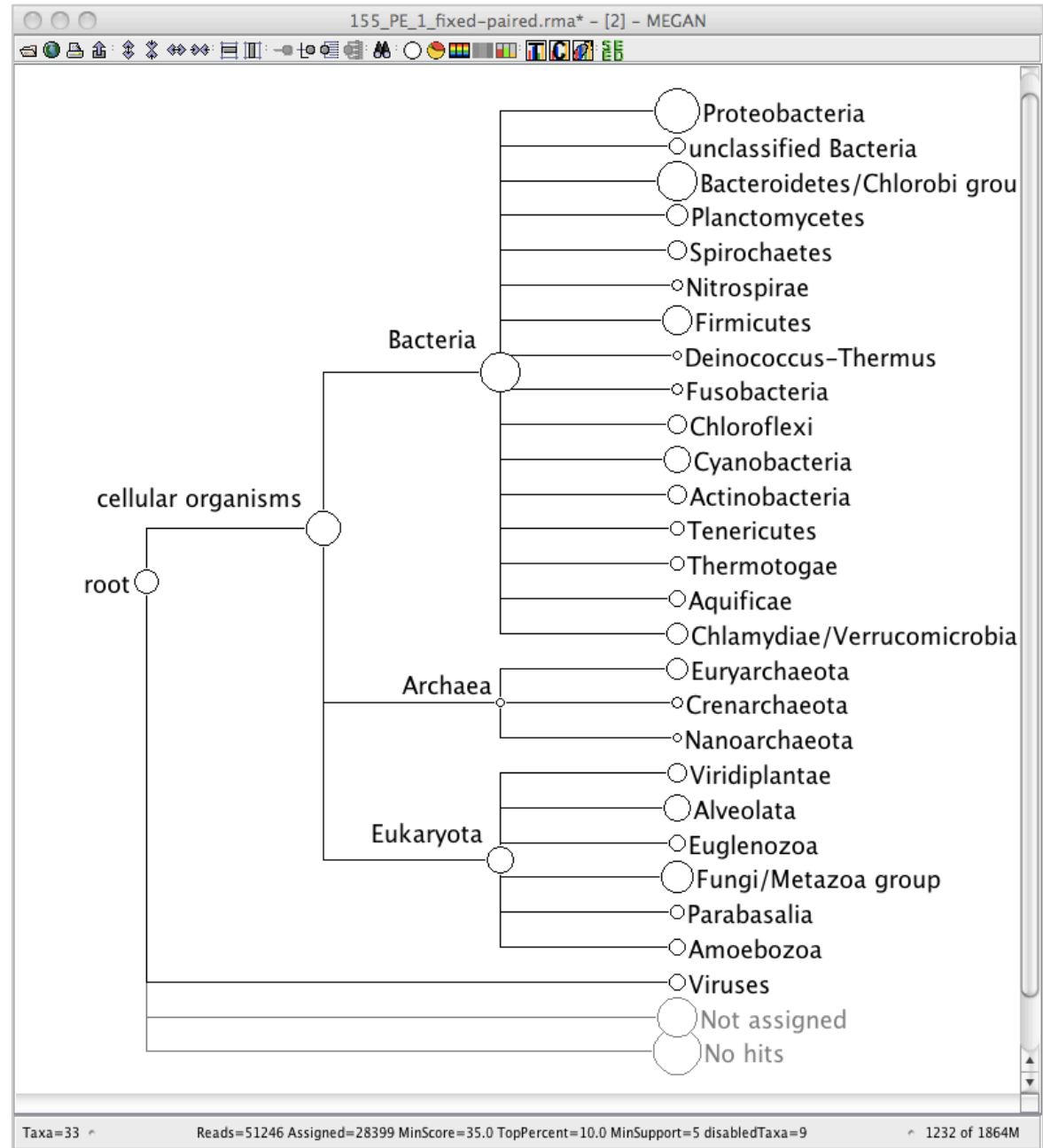
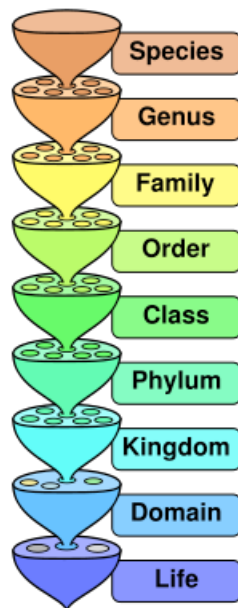
NCBI taxonomy:

- Contains all species represented by some sequence
- >560,000 nodes
- (2007: 280,000 nodes)
- King Phillip Came Over For Green Soup... (and more)



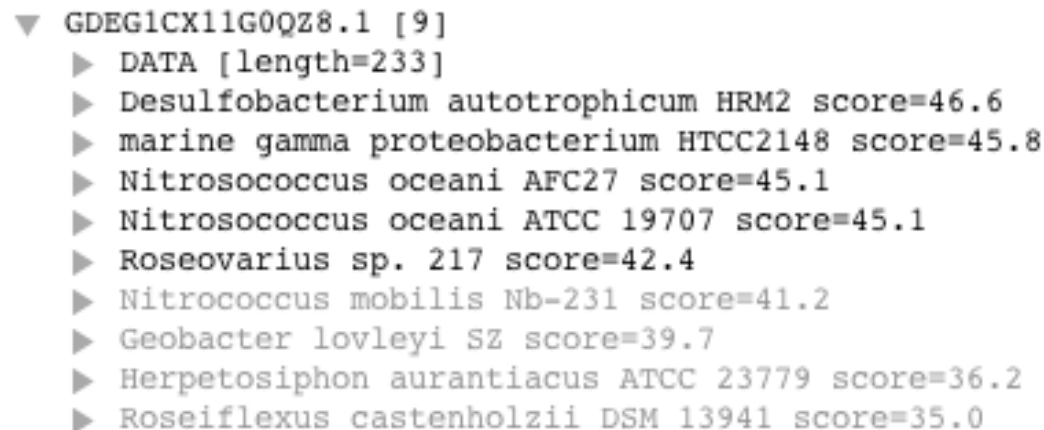
Organize and Visualize

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



Taxonomic analysis of 50,000 reads

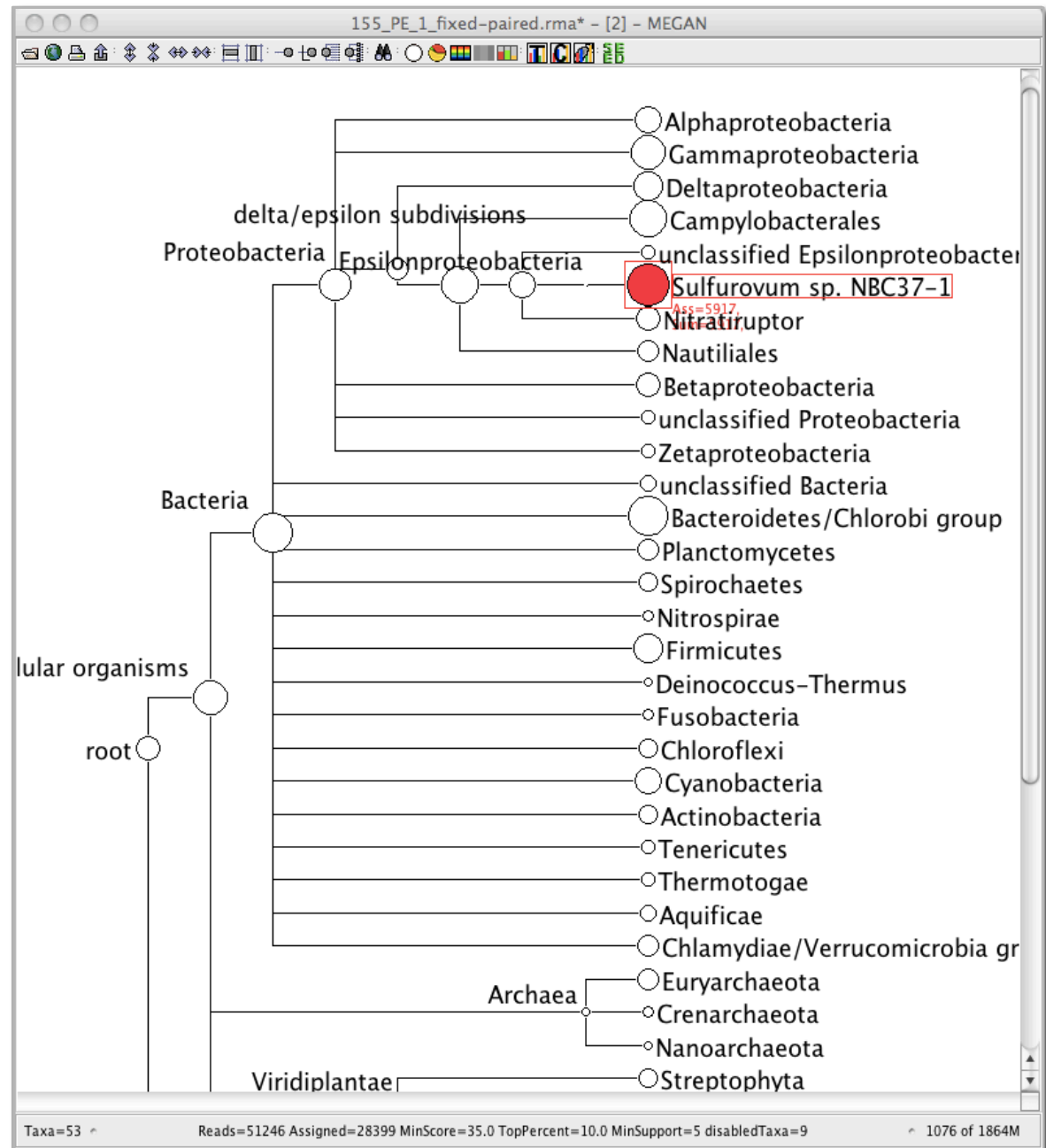
- ✓ Use NCBI taxonomy to bin sequences by evolutionary relatedness of organisms



Taxonomic analysis of 50,000 reads

Interact and Summarize

- ✓ Search for nodes of interest



Taxonomic analysis of 50,000 reads

Interact and Summarize

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

155_PE_1_fixed-paired.rma* - [2] - MEGAN

Inspector - 155_PE_1_fixed-paired.rma - [2] - MEGAN

155_PE_1_fixed-paired.rma

- ▼ Sulfurovum sp. NBC37-1 [5917]
 - ▶ GDEG1CX11G3K5K.2 [6]
 - ▶ GDEG1CX11G2QRJ.2 [100]
 - ▶ GDEG1CX11GZA6G.2 [22]
 - ▶ GDEG1CX11GTWQ6.1 [100]
 - ▶ GDEG1CX11G2DYH.1 [100]
 - ▶ GDEG1CX11G14WB.2 [1]
 - ▶ GDEG1CX11GZZGJ.1 [13]
 - ▶ GDEG1CX11G104U.1 [1]
 - ▶ GDEG1CX11GWIOM.2 [22]
 - ▶ GDEG1CX11GYO2L.2 [5]
 - ▶ GDEG1CX11G25HO.2 [100]
 - ▶ GDEG1CX11GY5Q9.1 [10]
 - ▶ GDEG1CX11G1VMN.2 [1]
 - ▶ GDEG1CX11GZR4E.1 [27]
 - ▶ GDEG1CX11GZR4E.2 [100]
 - ▶ GDEG1CX11G1W26.1 [100]
 - ▶ GDEG1CX11G1L5B.2 [4]
 - ▶ GDEG1CX11G1Y2D.1 [1]
 - ▶ GDEG1CX11G0Q6M.1 [100]
 - ▼ GDEG1CX11G2DJ2.1 [36]
 - ▼ DATA [length=237]
 - >GDEG1CX11G2DJ2.1 length=299 xy=2780 1424 region=11 run=R 2010 03 03 14 54 18
 - TTTtagagatagaaattatcccagacttaggagtagccatagacgaacttttcgctgcaatagacaaaagtataa
 - ▼ Sulfurovum sp. NBC37-1 score=80.5
 - >ref|YP_001357841.1| hypothetical protein SUN_0524 [Sulfurovum sp. NBC37-1]
 - dbj|BAF71484.1| conserved hypothetical protein [Sulfurovum sp. NBC37-1]
 - Length = 87
 - Score = 80.5 bits (197), Expect(2) = 8e-22
 - Identities = 36/51 (70%), Positives = 47/51 (92%)
 - Frame = +2
 - Query: 74 NATKEQEEDLEEMRDMRTECF+IIEELGRDELEAEEIDELLAELVEMKTEE 226
 - NATKEQ+EDLEEMR+MRTECF+I+EE+ +DEL+ EE +ELL ELV++KT+E
 - Sbjct: 35 NATKEQKEDLEEMRDMRTECF+IIEELGRDELEAEEIDELLAELVEMKTEE 85
 - Score = 46.6 bits (109), Expect(2) = 8e-22
 - Identities = 23/32 (71%), Positives = 27/32 (84%)
 - Frame = +3
 - Query: 3 LEIEIIPDLGVAIDELFAAIDKSKMPQKNKKK 98
 - LE+EIIPDL VAIDELFAAIDK+K K +K+
 - Sbjct: 11 LEIEIIPDLGVAIDELFAAIDKAKNATKEQKE 42
 - ▶ Strongylocentrotus purpuratus score=37.0
 - ▶ Drosophila erecta score=35.8
 - ▶ Rattus norvegicus score=35.8
 - ▶ Algoriphagus sp. PR1 score=35.4
 - ▶ Gallus gallus score=35.0

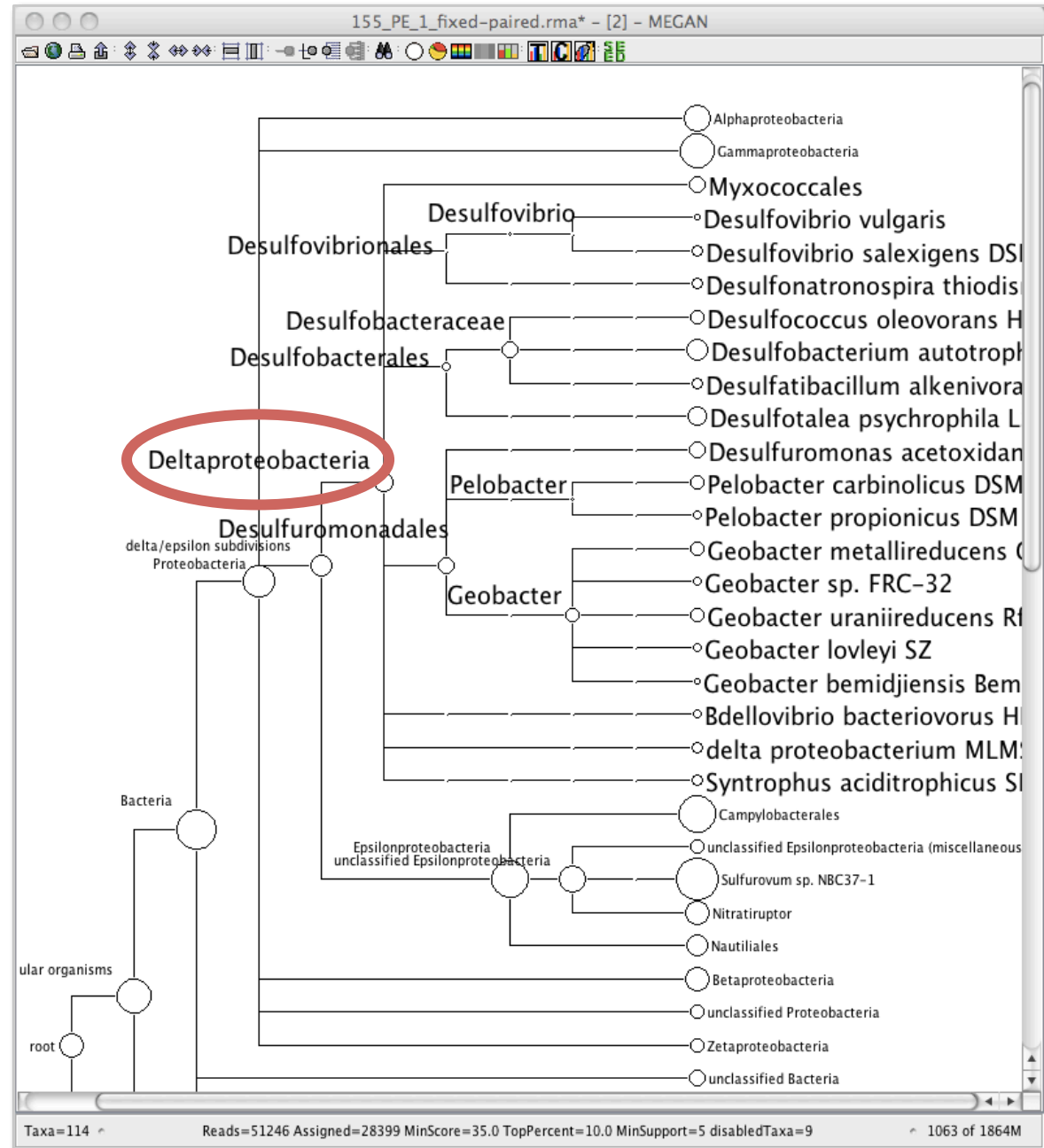
Viridiplantae Streptophyta

Taxa=53 Reads=51246 Assigned=28399 MinScore=35.0 TopPercent=10.0 MinSupport=5 disabledTaxa=9 1076 of 1864M

Taxonomic analysis of 50,000 reads

Interact and Summarize

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

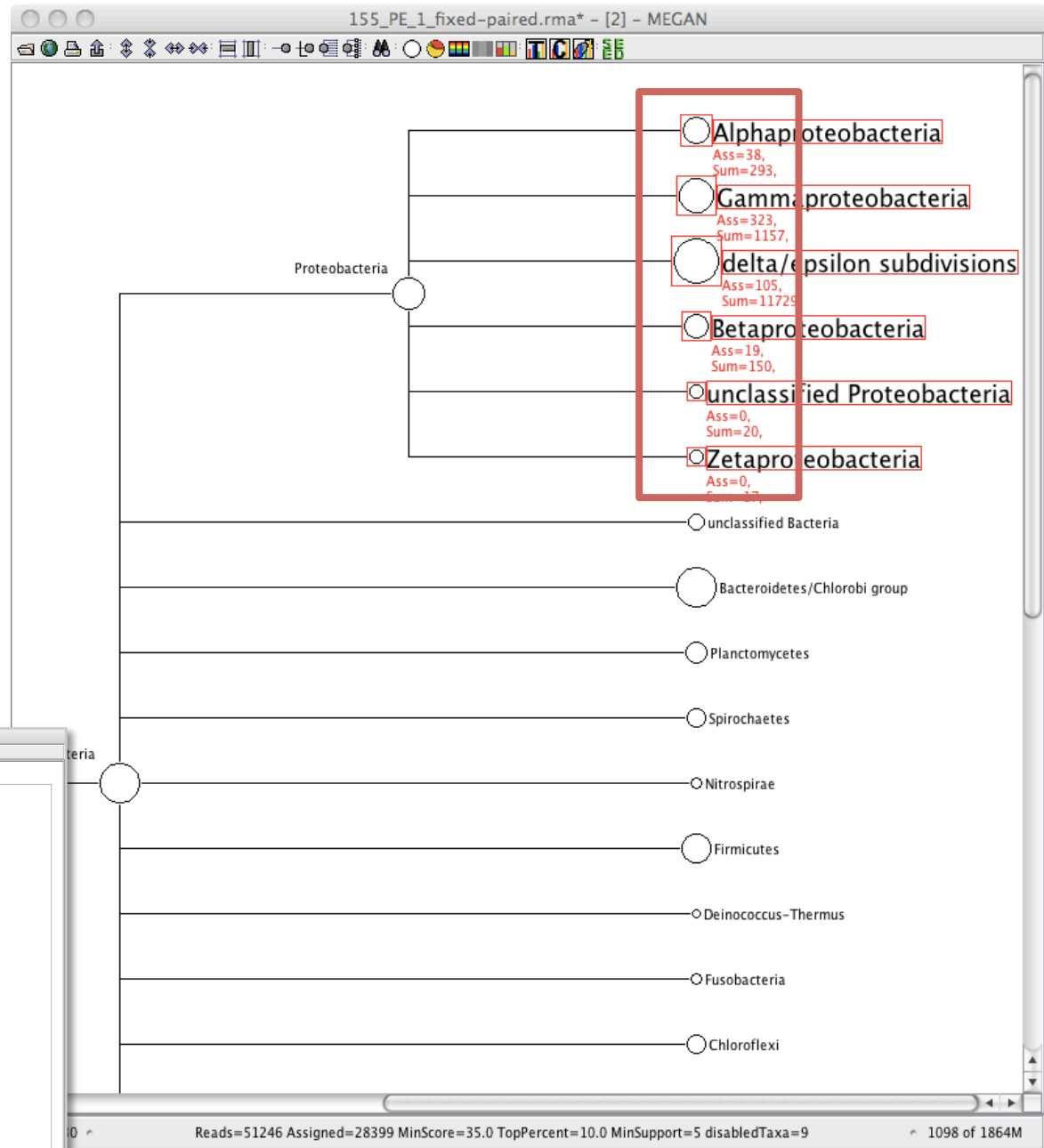
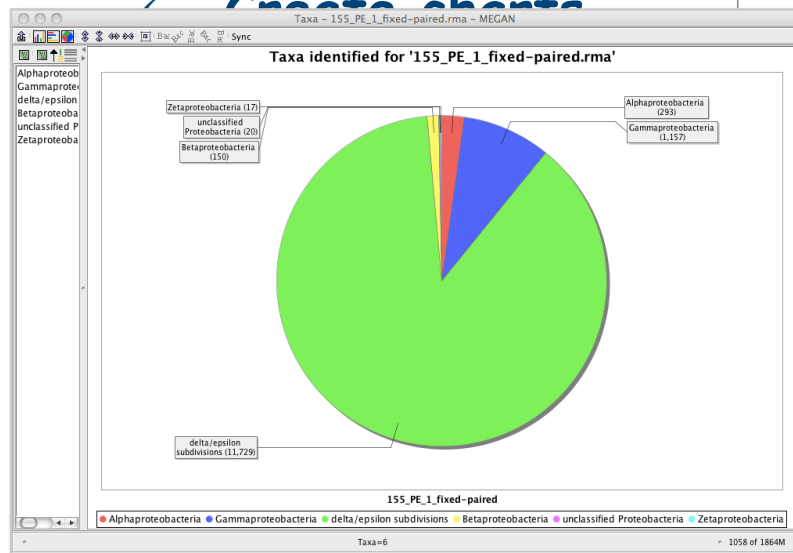


Taxonomic analysis of 50,000 reads

Interact and Summarize

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

Create charts



Taxonomic analysis of 50,000 reads

Capture

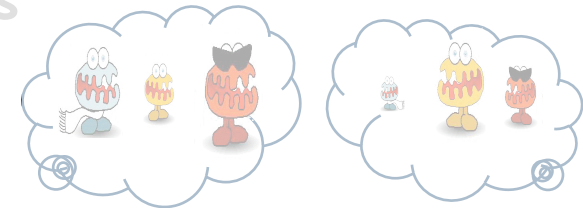
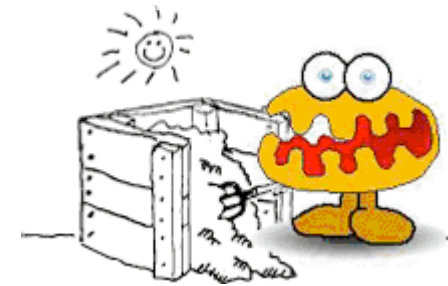
- ✓ Capture all sequences (and/or their matches) assigned to selected nodes

The screenshot displays the MEGAN software interface. The main window shows a taxonomic tree on the left and a list of sequences on the right. The 'Export' menu is open, showing options like 'Export', 'Export Image...', 'Page Setup...', 'Print...', 'Compare...', 'Extract Reads By Taxa...', 'Extract Reads By COGs...', 'Extract Reads By Subsystems...', 'Import CSV...', 'Tools', 'Properties...', and 'Close'. The 'Extract Reads By Taxa...' option is highlighted. Below the main window, the 'Extract by Taxa - 155_PE_1_fixed-paired.rma - [2] - MEGAN' dialog box is open. It shows the 'Output files' section with 'Directory: /Users/huson/data/megan' and 'File name: reads-%t.fasta'. The 'Include Summarized' checkbox is checked. The 'Execute' button is visible. The status bar at the bottom of the MEGAN window shows 'Taxa=1204', 'Reads=51246', 'Assigned=28399', 'MinScore=35.0', 'TopPercent=10.0', 'MinSupport=5', 'disabledTaxa=9', and '987 of 1864M'.

Taxonomic analysis of 50,000 reads

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



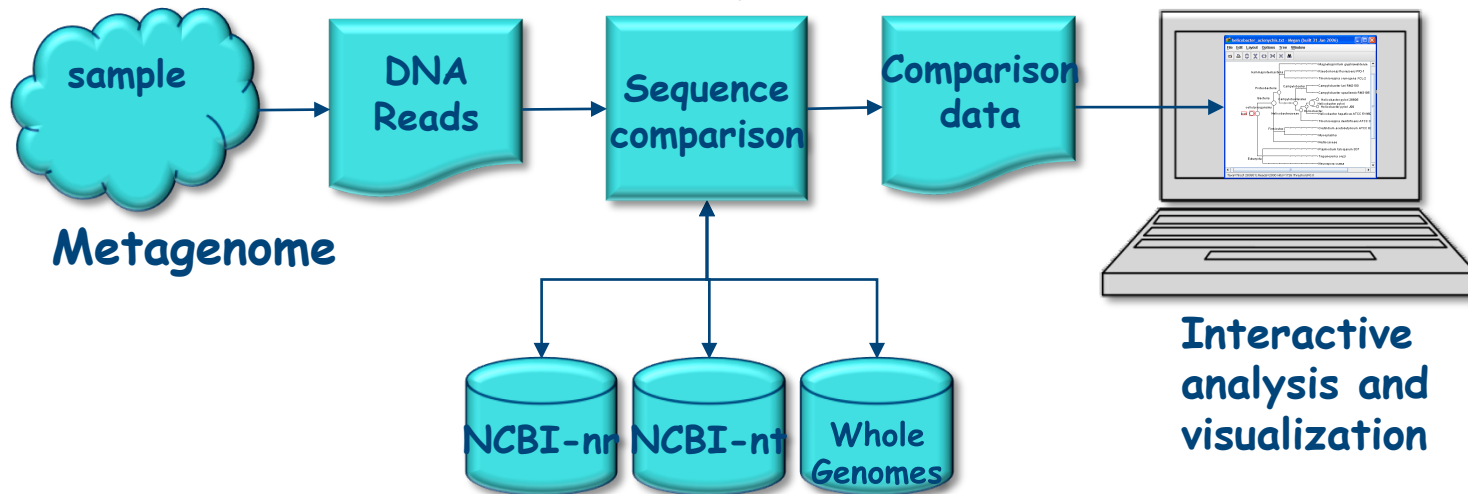


Identifying Taxa and Genes

Metagenome analysis

Basic idea: compare reads against references sequences of known species and/or function

**BLASTX
against
NCBI-NR**



Reference databases

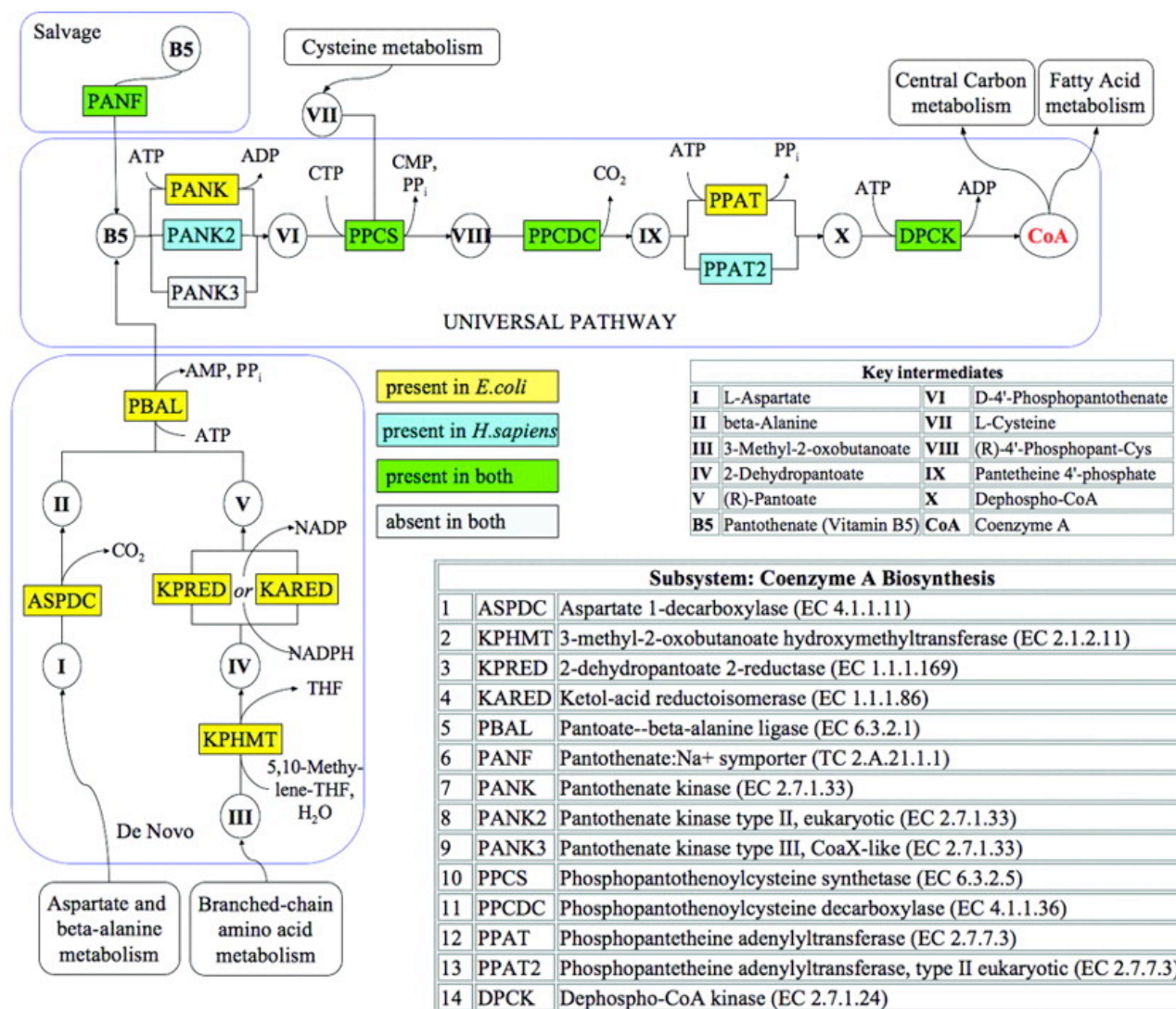


Functional Analysis using SEED

- The SEED classification assigns genes to functional roles in subsystems
- A subsystem is a set of functional roles that implement a specific biological process or structural complex
- RAST and MG-RAST: Rapid annotation using subsystem technology
- Graph has ~10,000 nodes and edges
- www.theSEED.org

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

Example of a Subsystem



Coenzyme A Biosynthesis Subsystem

Organize and Visualize

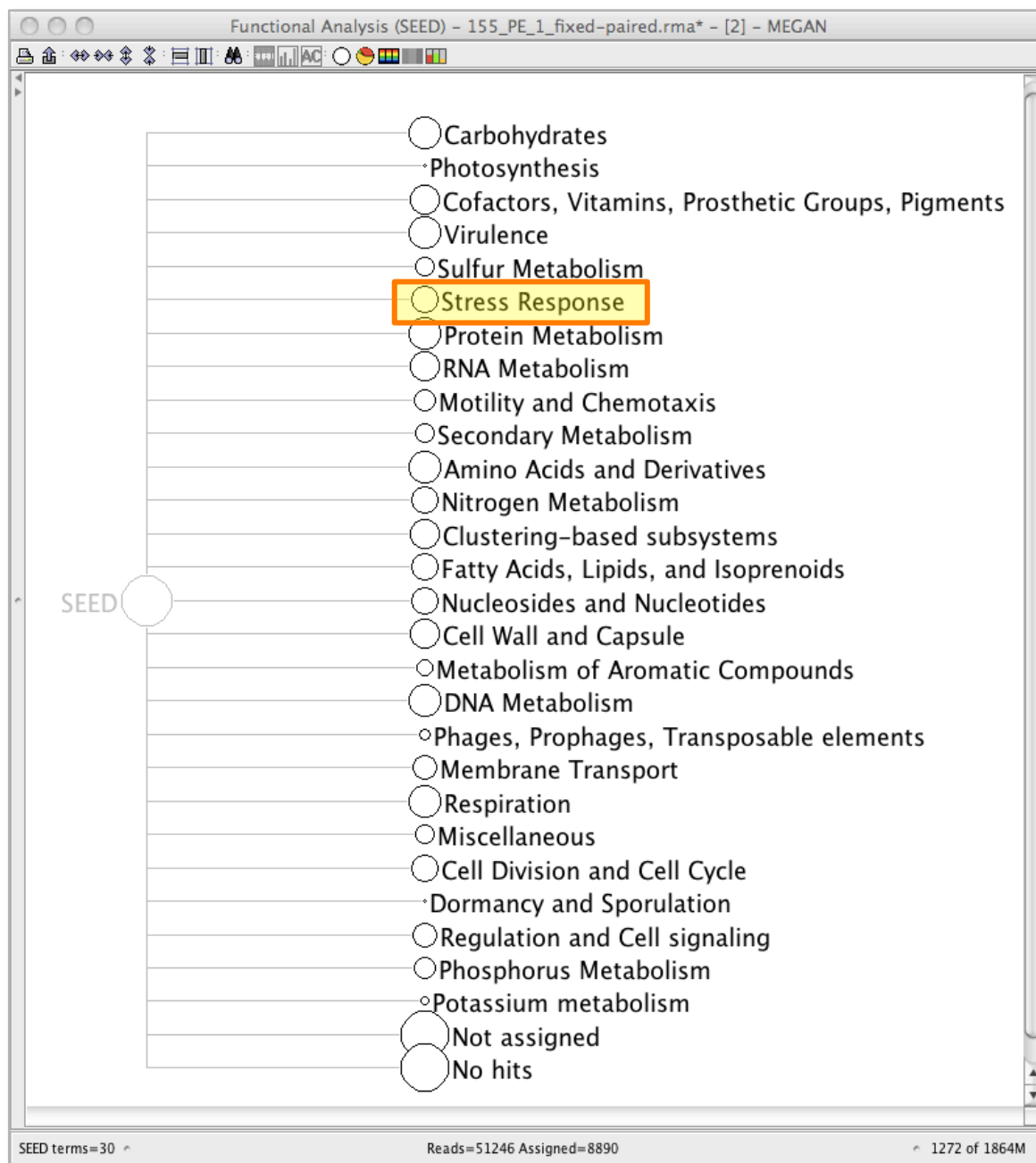
Functional analysis

✓ Use SEED
classification to bin
sequences by
subsystems



www.theseed.org

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



SEED analysis of 50,000 reads

Organize and Visualize

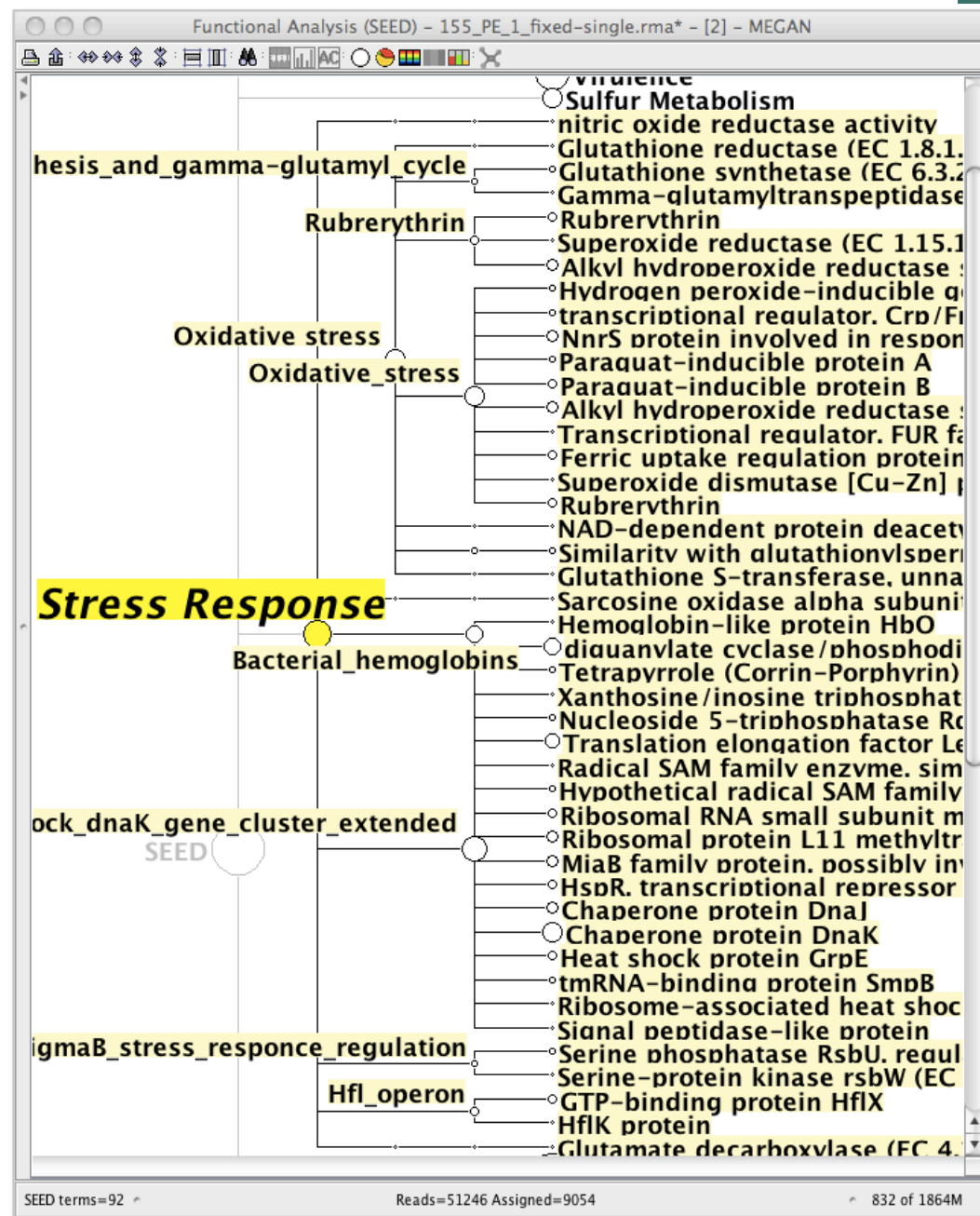
Functional analysis

- ✓ Use SEED classification to bin sequences by subsystems
- ✓ ... and by functional roles



www.theseed.org

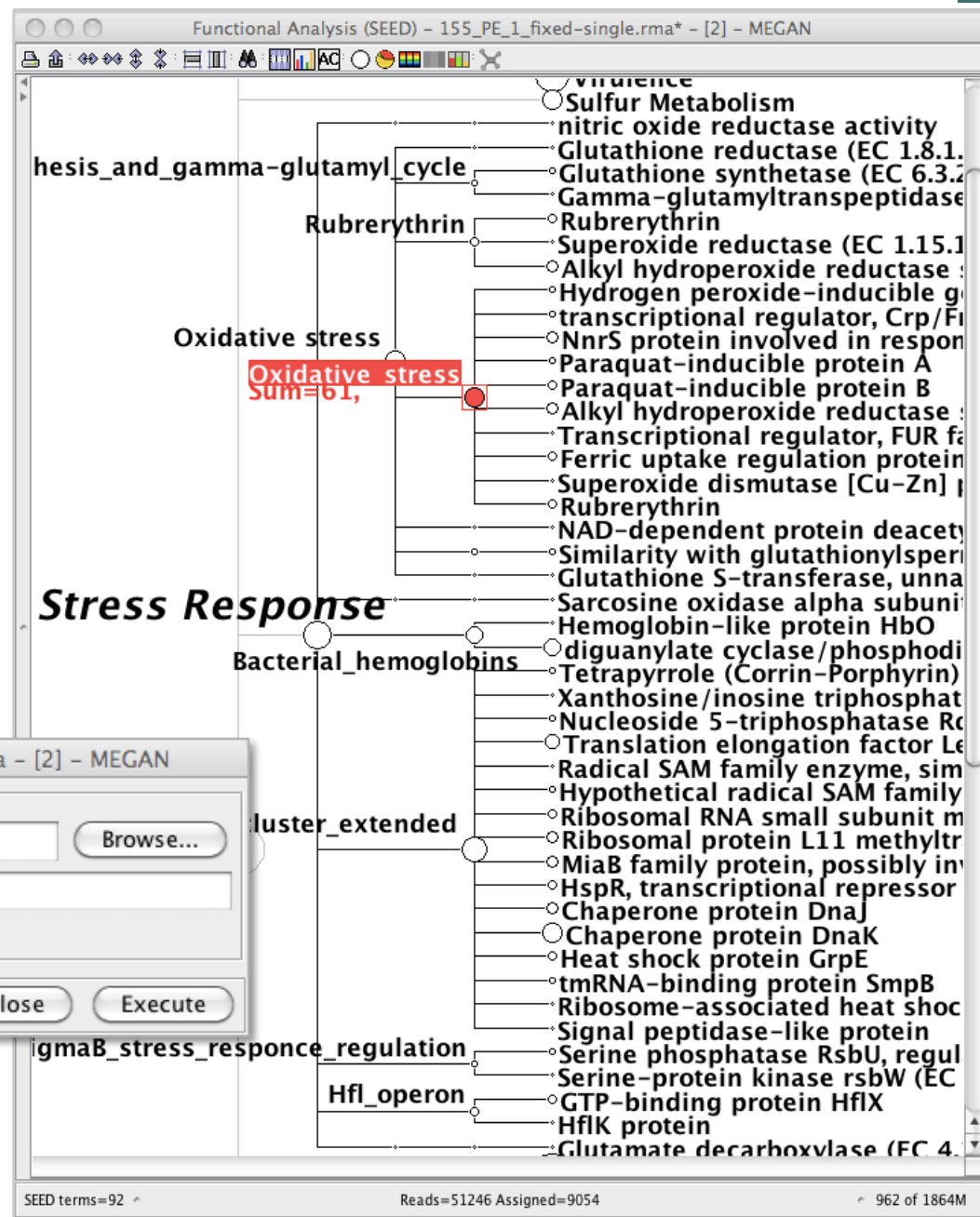
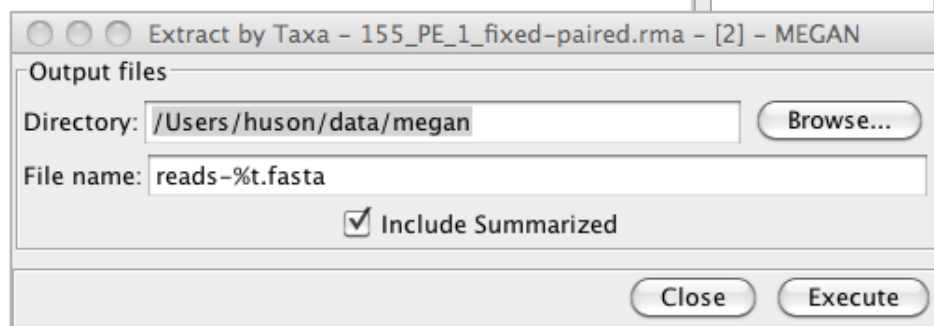
SEED: Overbeek et al.,
Nucleic Acids Res 33 (17), 2005



SEED analysis of 50,000 reads

Capture

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



SEED analysis of 50,000 reads

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



MEGAN - KEGG Viewer

Statistics T1B1 T1B6 T2B6 T2B1 Overview

STARCH AND SUCROSE METABOLISM

00500 6/18/10
(c) Kanehisa Laboratories

MEGAN KEGG Viewer

map

71 of 992M

Daniel Huson © 2010



Three Basic Computational Questions

- Who is out there?

- Types of organisms
- In what proportions?



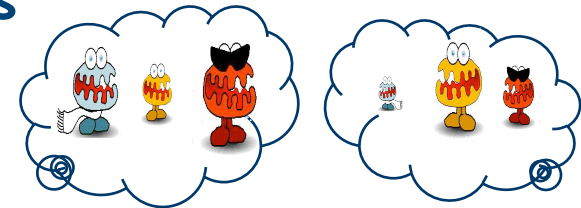
- What are they doing?

- Types of genes
- Which metabolic pathways?
- In what proportions?



- How do different samples compare?

- Pairwise and multiple comparisons
- Correlations with environmental parameters?

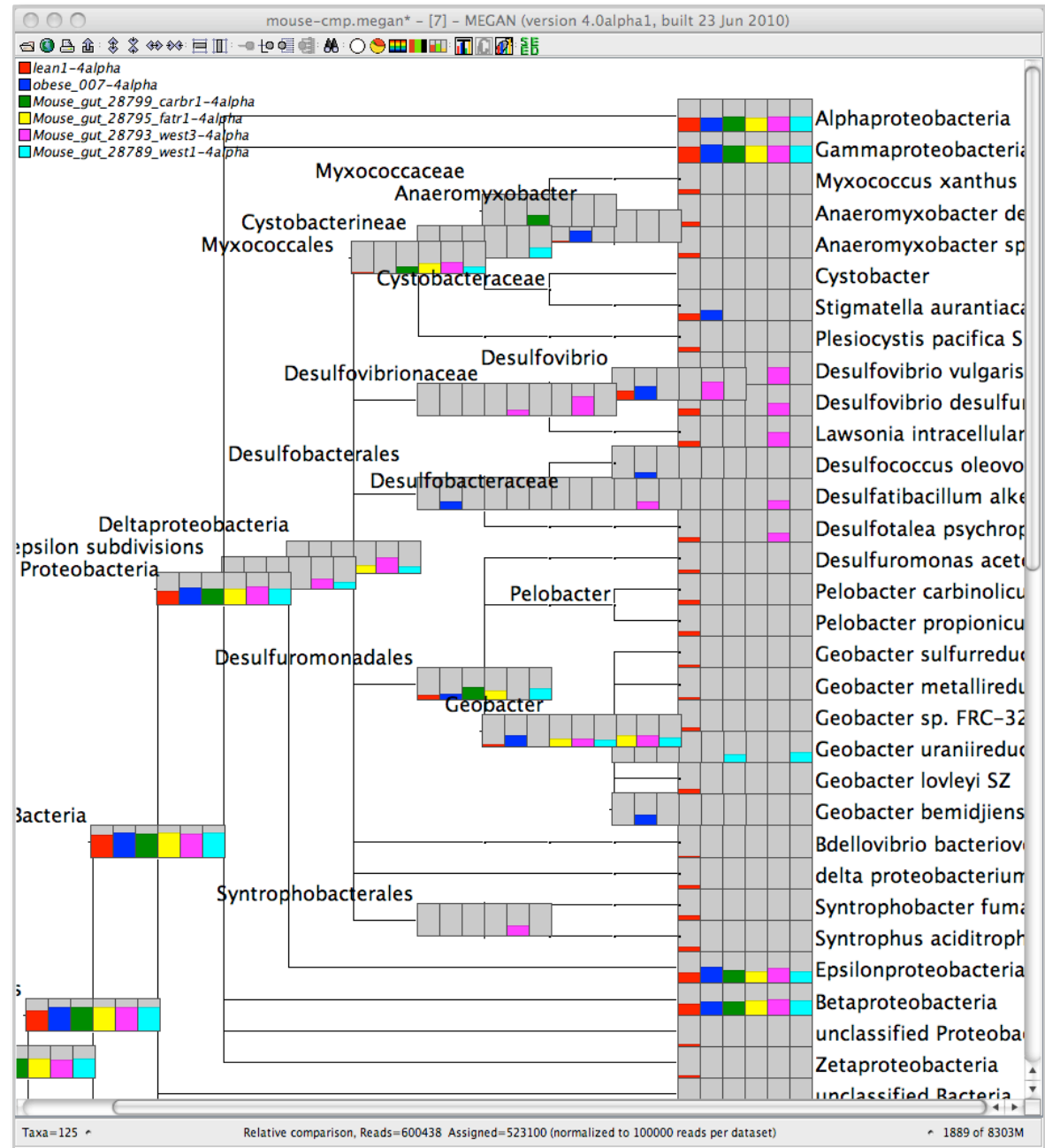
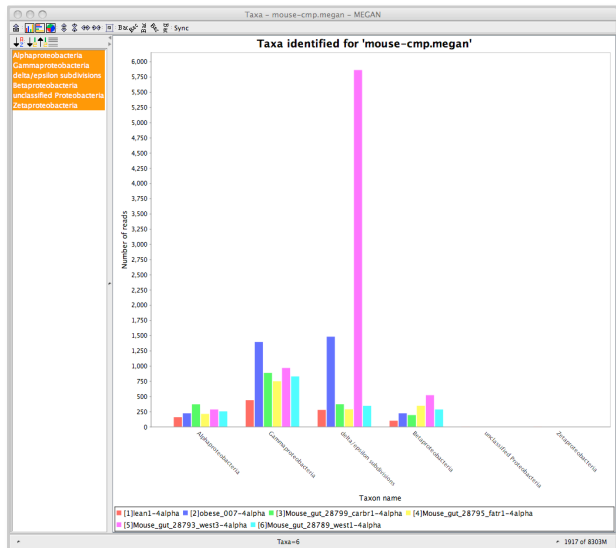


- Serve to answer biological or medical questions

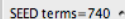
Compare

Display multiple
datasets
simultaneously

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



✓ ... and summarize



Compare

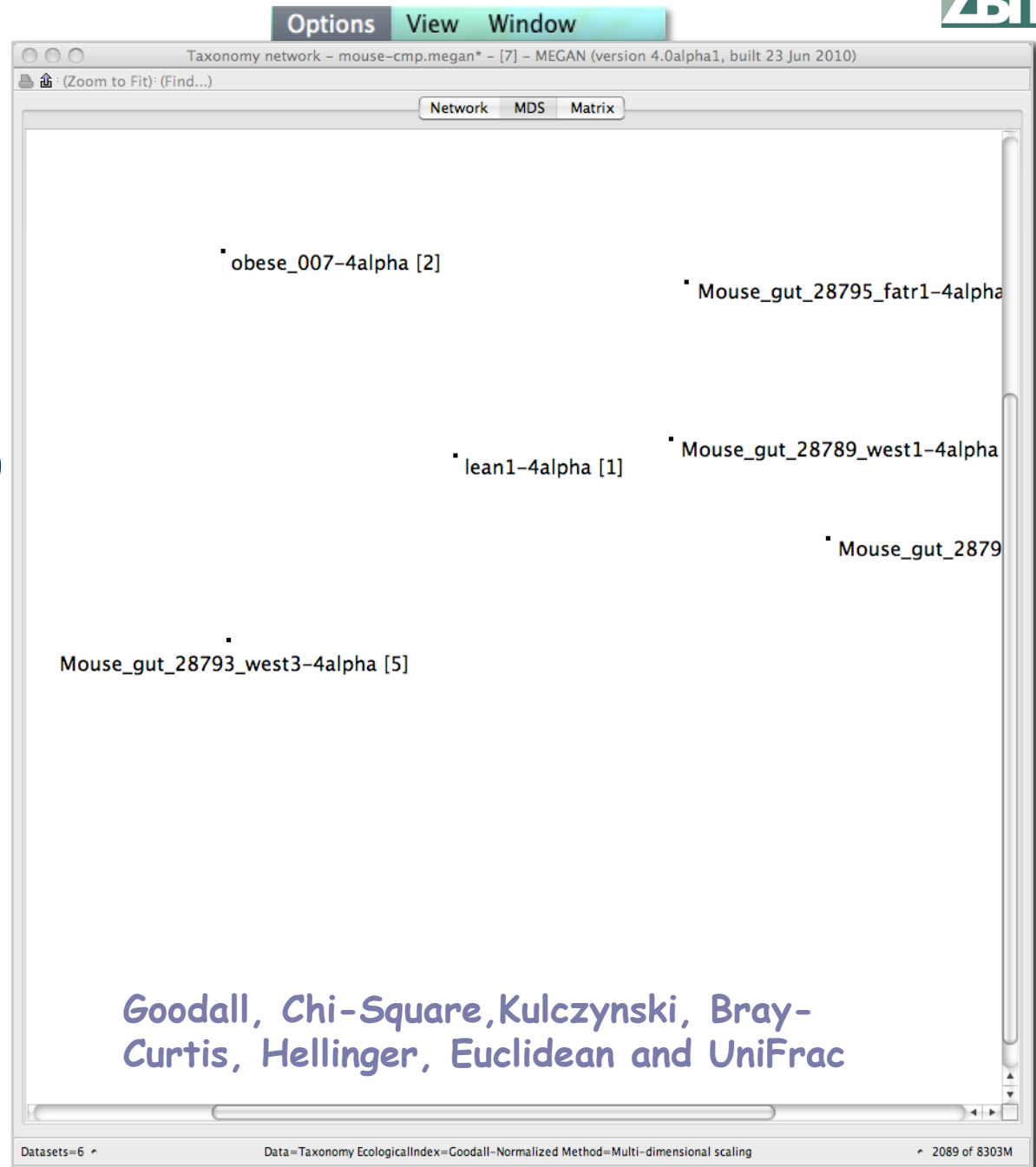
High-level comparison:

- ✓ 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

Daniel Huson © 2010

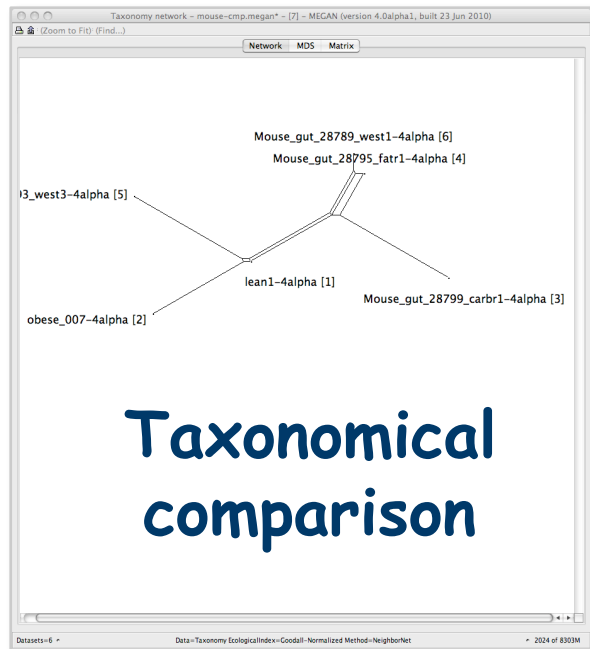


Goodall, Chi-Square, Kulczynski, Bray-Curtis, Hellinger, Euclidean and UniFrac

Compare

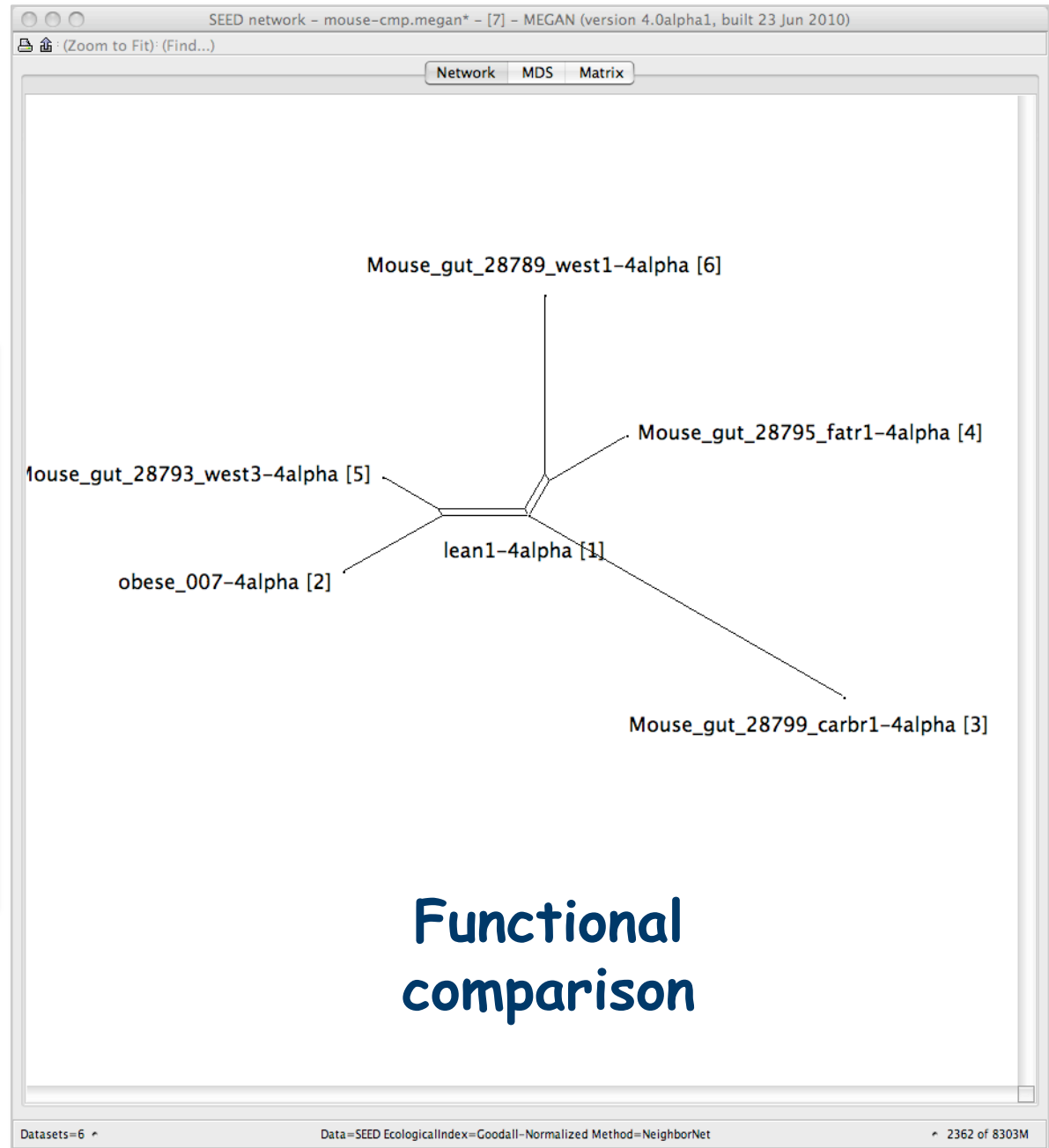
High-level comparison:

✓ Select taxa



**Taxonomical
comparison**

✓ Select functions...

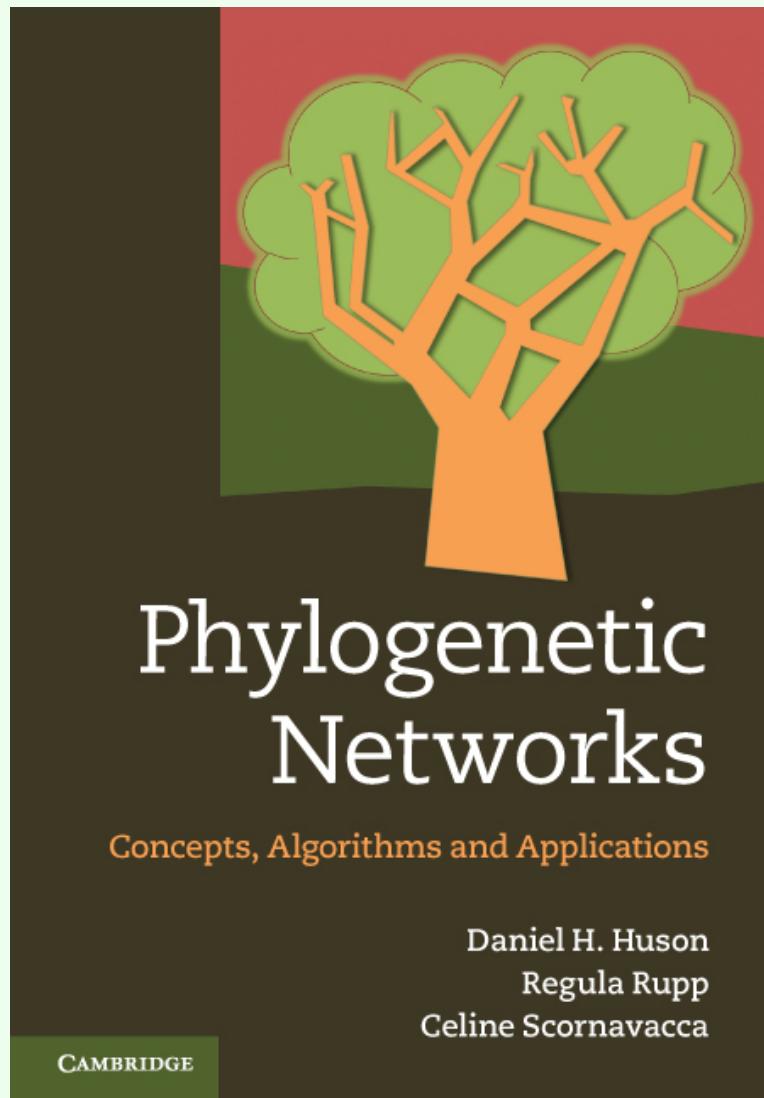


**Functional
comparison**



Phylogenetic Networks

Concepts, Algorithms and Applications



~ 360 pages

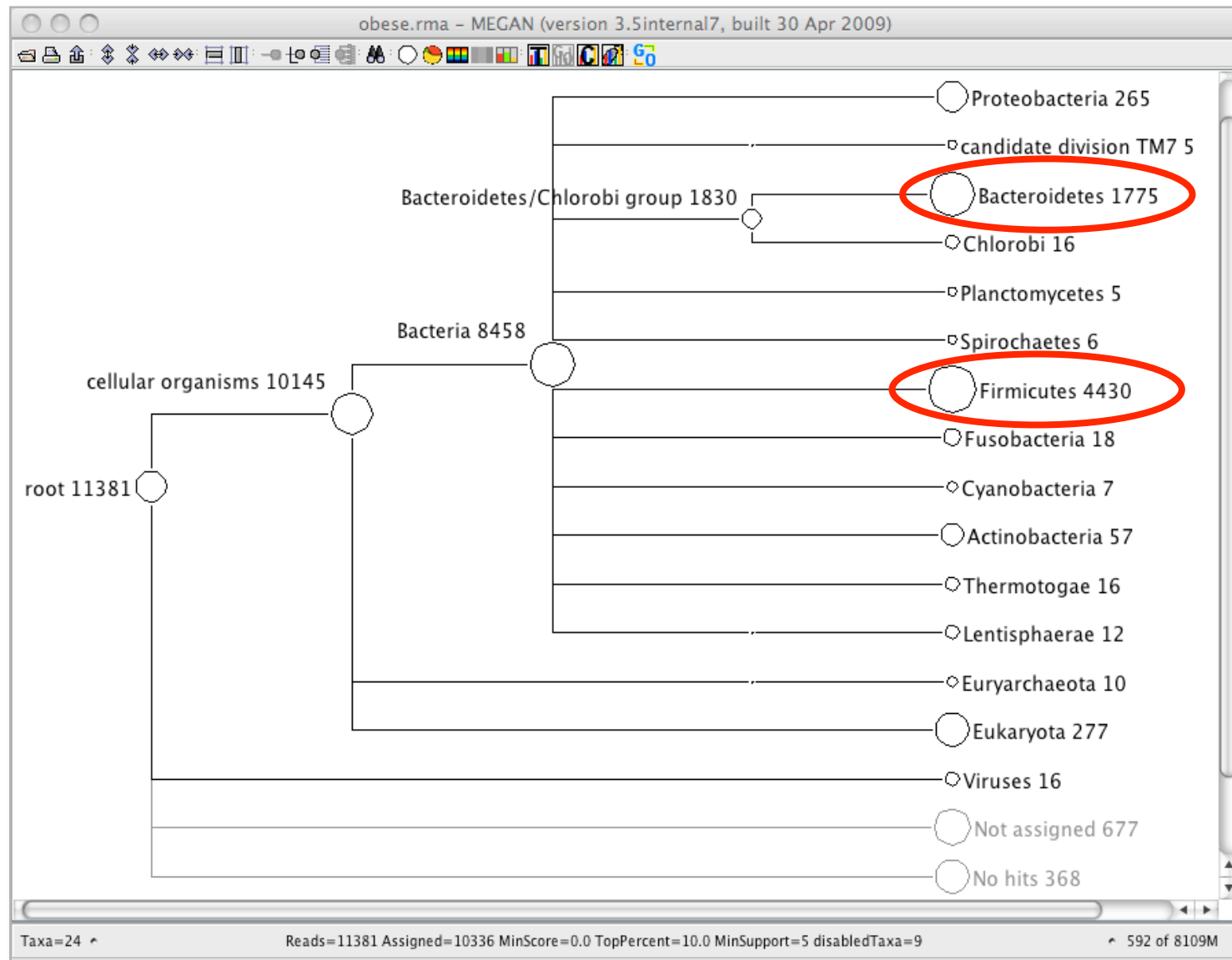
- 55 lemmas
- 20 theorems
- 50 algorithms
- 85 exercises
- 15 applications
- 190 figures

~ 40 EUR

Dec 2010

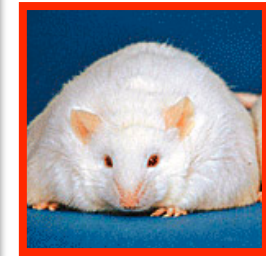
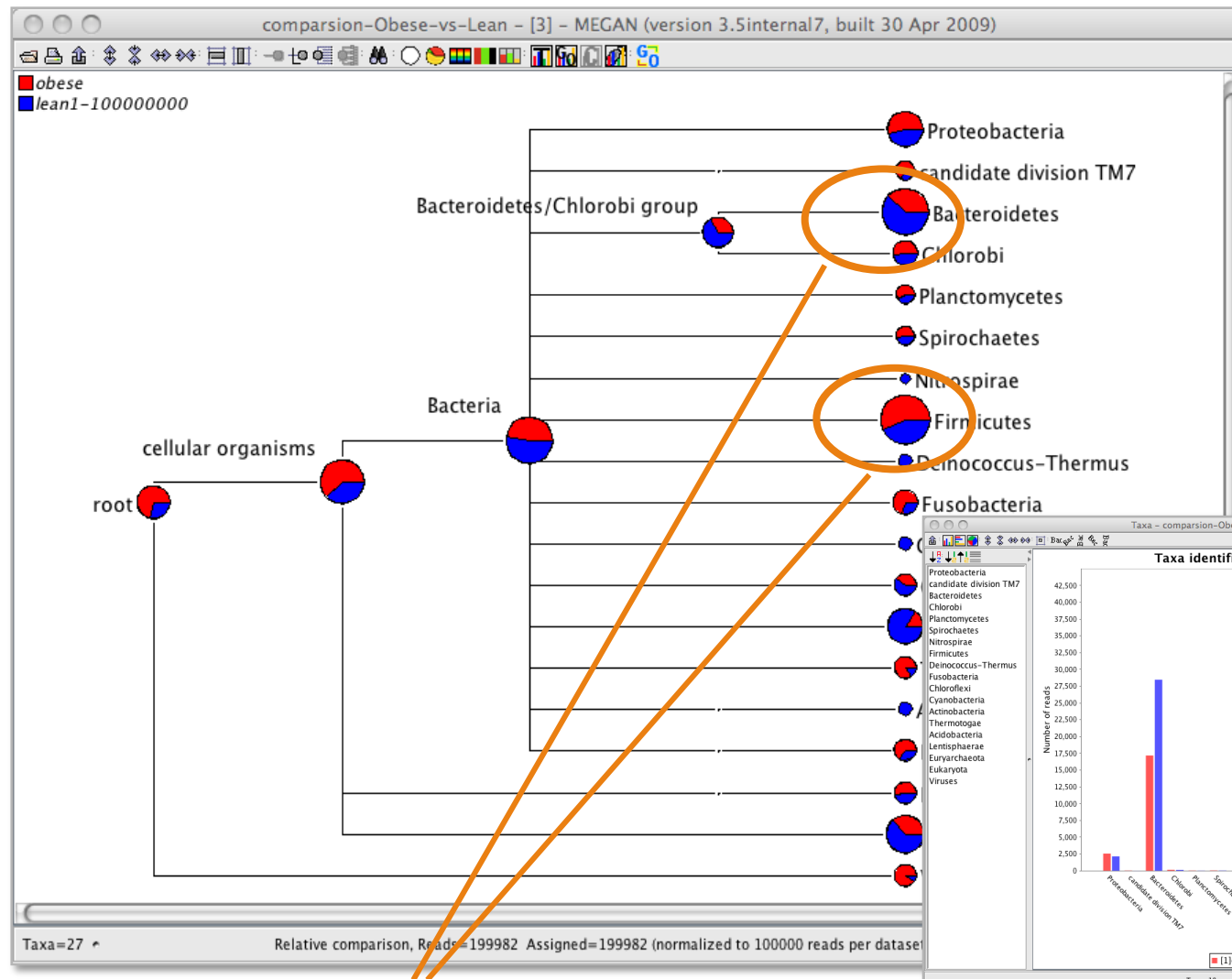


Example: Mouse Gut Microbiome

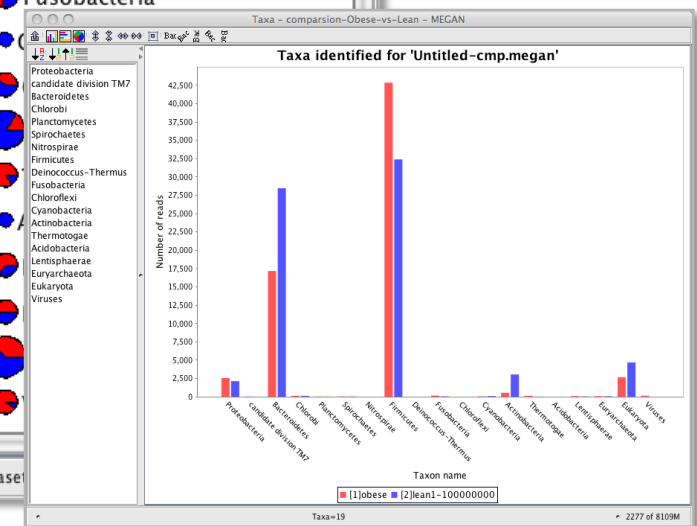


Dominant gut microbes: Bacteroidetes and Firmicutes

Comparative Analysis



Change in proportions of these two phyla





Basic Principles

How to
analyze
a
meta
genome?

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



Contents

- Genomics
- Sequencing
- Metagenomics
- **More computational questions**
- Outlook



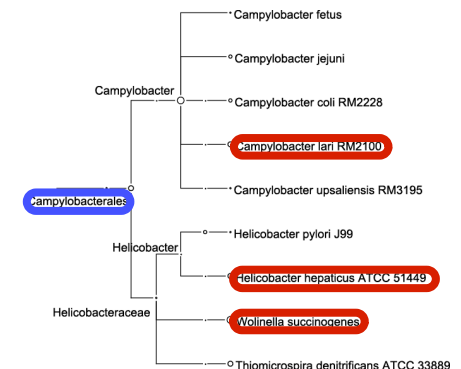
Single Reads vs Paired Reads

In metagenomics:

- Use single reads or paired reads?
- In the latter case, short clones or long clones?

Taxonomical Analysis Based on Gene Content

Set of
all species
containing gene



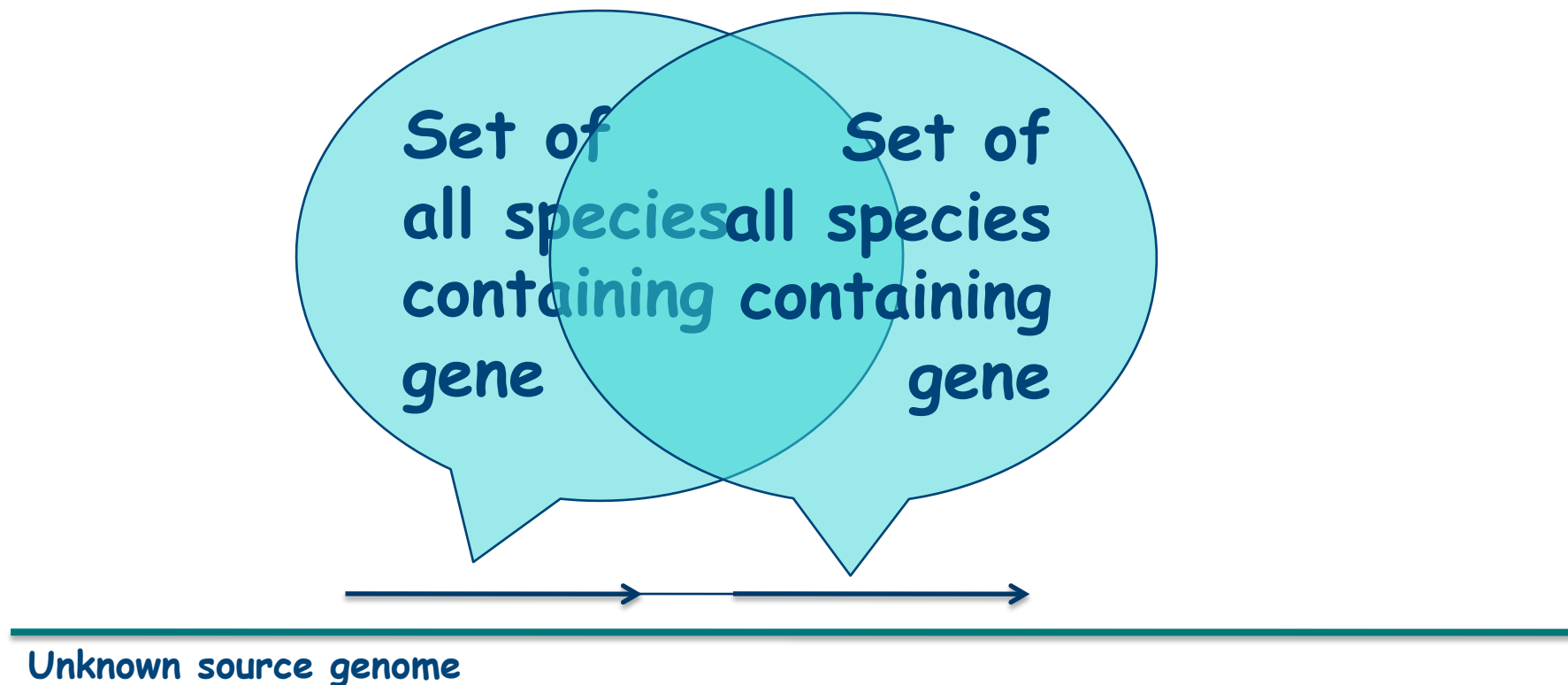
read

Unknown source genome

Use LCA to assign to (higher-rank) taxon

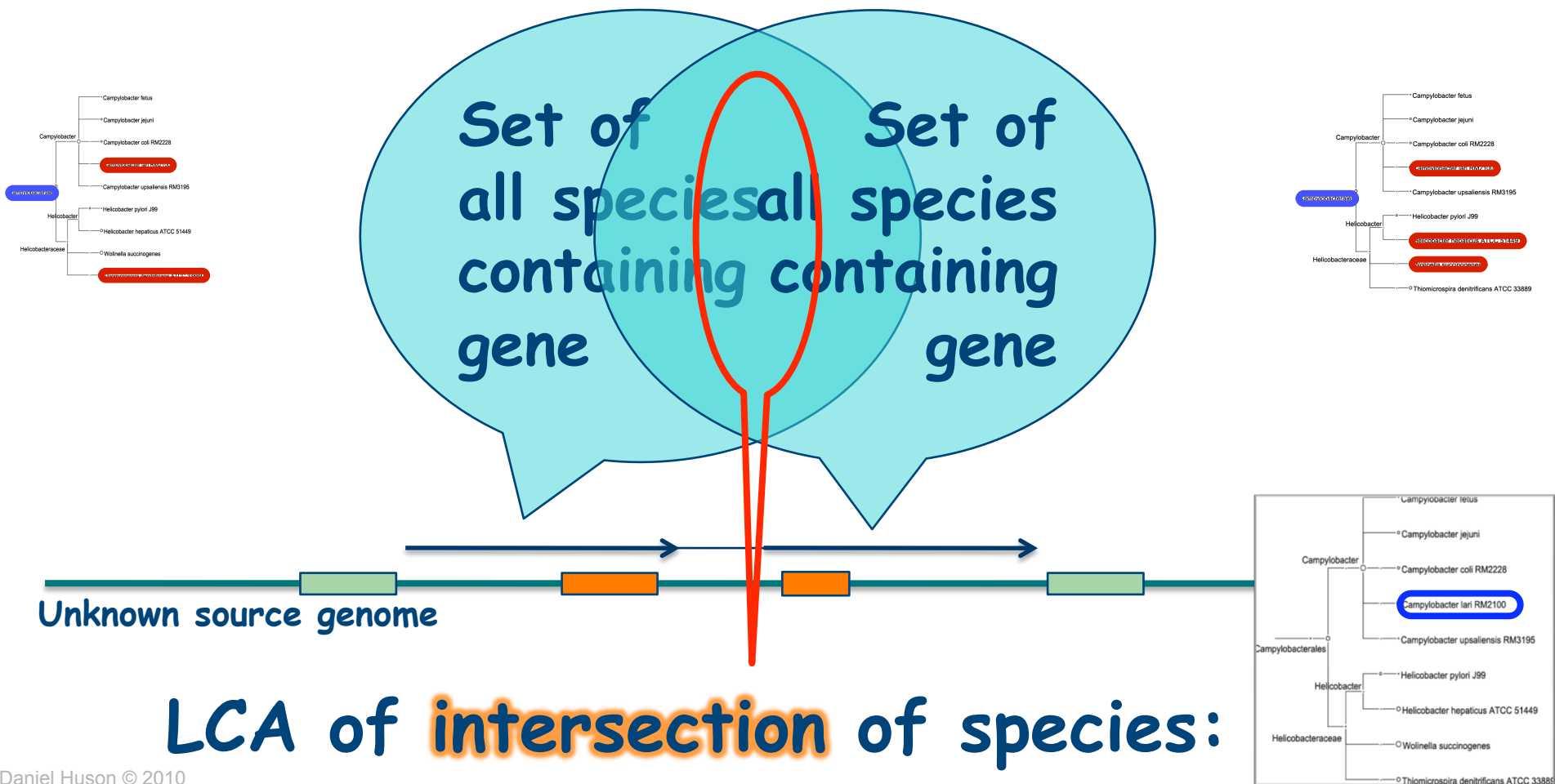
Taxonomical Analysis Based on Gene Content

Short clones or long clones?

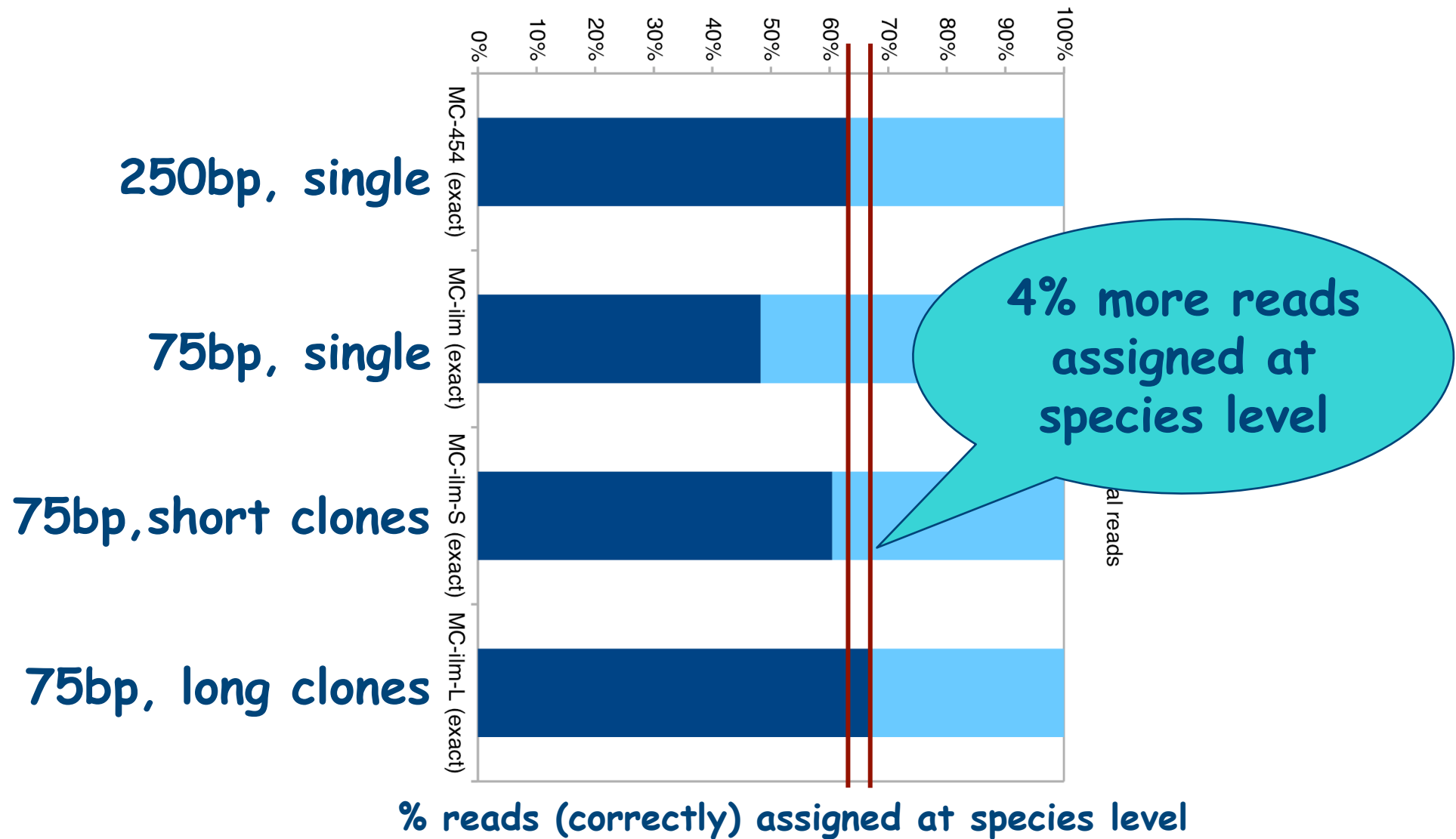


Taxonomical Analysis Based on Gene Content

Claim: long clones are more specific



Simulated Performance 454 vs Illumina

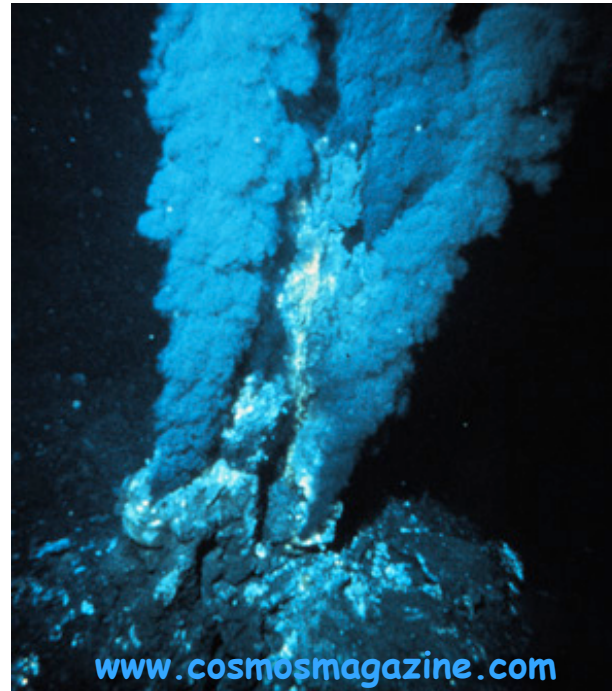


MetaSim - simulator (Richter et al, 2008)



Ocean Seabed Sample

- Joint work with Ida Steen (Bergen)



- Currently analyzing large set of 454 paired reads (7kb clones) collected from an Arctic hot vent at 3km depth



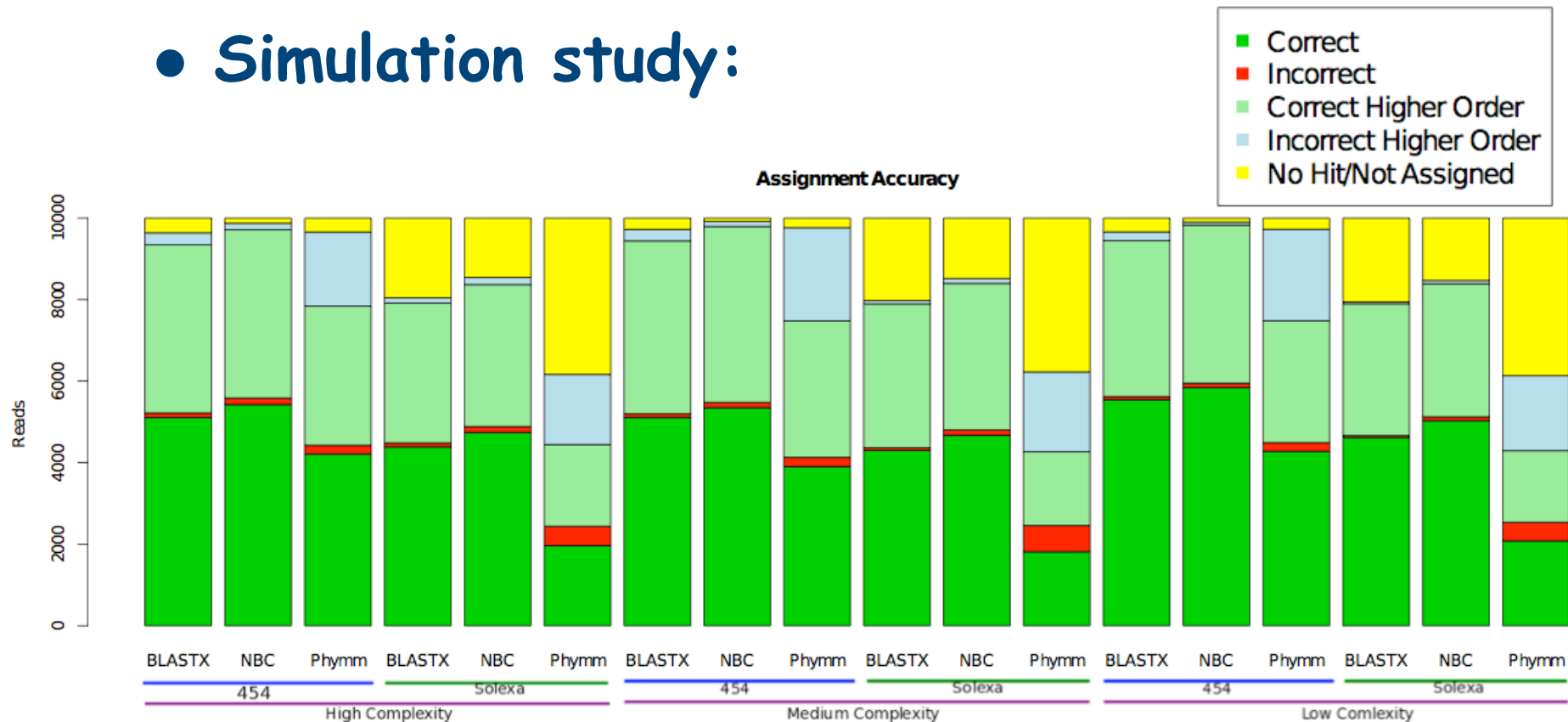
Hybrid Approach

- Machine-learning based classification approaches are much faster than BLASTX
- Biologists want to see alignments and they are needed for functional analysis
- Hybrid approach:
 - Use taxonomic classifier to perform taxonomic binning
 - BLASTX reads only against assigned taxa
- Study of NBC (Rosen *et al* 2008) and Phymm (Brady & Salzberg, 2009)



Hybrid Approach

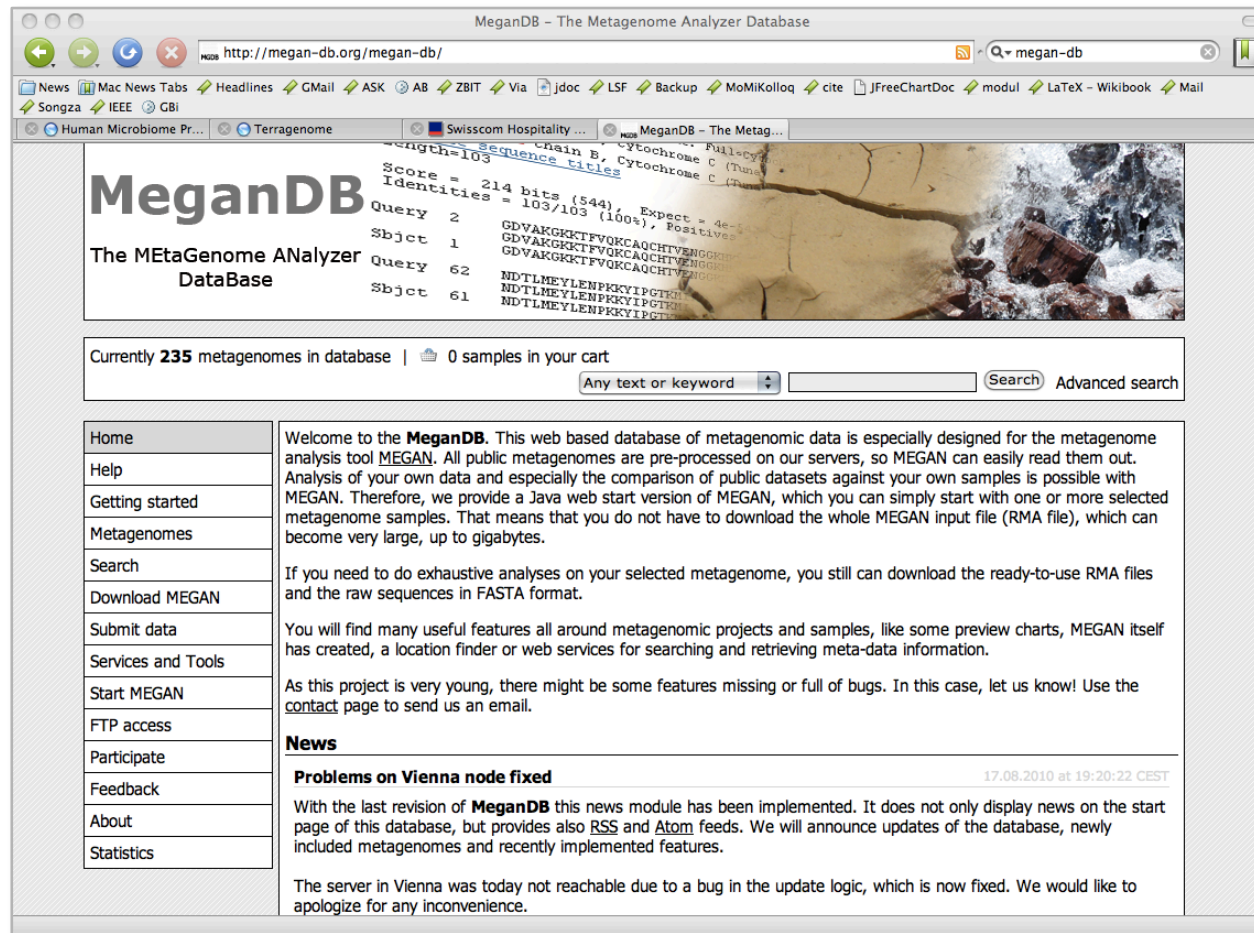
- Simulation study:



- ~10x speed-up over full BLASTX
- Increased accuracy using NBC
- Decreased accuracy using Phymm

MEGAN-DB

- Database of precomputed MEGAN files



- www.megan-db.org

Joint work with Thomas Rattei and Simon Domke



Contents

- Genomics
- Sequencing
- Metagenomics
- Computational questions
- **Outlook**

Computational Challenges...

- **Global Ocean Sampling**
 - www.jcvi.org/cms/research/projects/gos
- **Human Microbiome Project**
 - nihroadmap.nih.gov/hmp/
- **Terragenome Consortium**

Terabases of sequences



- **Survey of the Earth microbiome**

- Petabases of sequences
 - Processing and storage of exabytes of data
- (mega, giga, tera, peta, exa...)





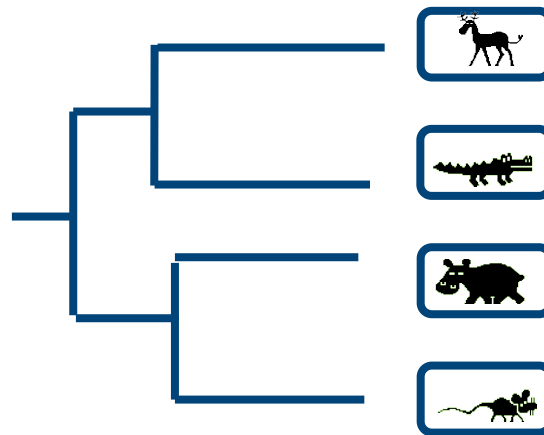
Computational Challenges...

- Data storage and access
- Tools for navigating metagenome data and metadata
- Ever faster analysis methods
- How to learn across multiple datasets?
- How to build a model of the Earth microbiome?



Joint Work With:

- **Tübingen:** Suparna Mitra, Daniel Richter, Nico Weber & Max Schubach
- **Penn State:** Stephan Schuster and Qi Ji
- **Vienna:** Tim Urich, Christa Schleper, Thomas Rattei, Simon Domke
- **Bergen:** Ida Steen and Anders Lanzen



www-ab.informatik.uni-tuebingen.de