

# Towards understanding the functional and taxonomic repertoire of microbial communities using the EBI metagenomics portal

Rob Finn ([rdf@ebi.ac.uk](mailto:rdf@ebi.ac.uk)),  
Open Door Workshop, Hinxton  
13th May 2015



# Metagenomics - a broad range of applications



# Metagenomics - a broad range of applications



CMAJ

RESEARCH

## Gut microbiota of healthy Canadian infants: profiles by mode of delivery and infant diet at 4 months

Meghan B. Azad PhD, Theodore Konya MPH, Heather Maughan PhD, David S. Guttman PhD, Catherine J. Field PhD, Radha S. Chari MD, Malcolm R. Sears MB, Allan B. Becker MD, James A. Scott PhD, Anita L. Kozyrskyj PhD, on behalf of the CHILD Study Investigators

## ARTICLE

doi:10.1038/nature11234

## Structure, function and diversity of the healthy human microbiome

The Human Microbiome Project Consortium\*

OPEN ACCESS Freely available online

PLOS PATHOGENS

## Targeted Restoration of the Intestinal Microbiota with a Simple, Defined Bacteriotherapy Resolves Relapsing *Clostridium difficile* Disease in Mice

Trevor D. Lawley<sup>1\*</sup>, Simon Clare<sup>1,3</sup>, Alan W. Walker<sup>1,3</sup>, Mark D. Stares<sup>1</sup>, Thomas R. Connor<sup>1</sup>, Claire Raisen<sup>1</sup>, David Goulding<sup>1</sup>, Roland Rad<sup>1</sup>, Fernanda Schreiber<sup>1</sup>, Cordelia Brandt<sup>1</sup>, Laura J. Deakin<sup>1</sup>, Derek J. Pickard<sup>1</sup>, Sylvia H. Duncan<sup>2</sup>, Harry J. Flint<sup>2</sup>, Taane G. Clark<sup>3</sup>, Julian Parkhill<sup>1</sup>, Gordon Dougan<sup>1</sup>

## LETTER

doi:10.1038/nature11582

## Host-microbe interactions have shaped the genetic architecture of inflammatory bowel disease

OPEN ACCESS Freely available online

PLOS ONE

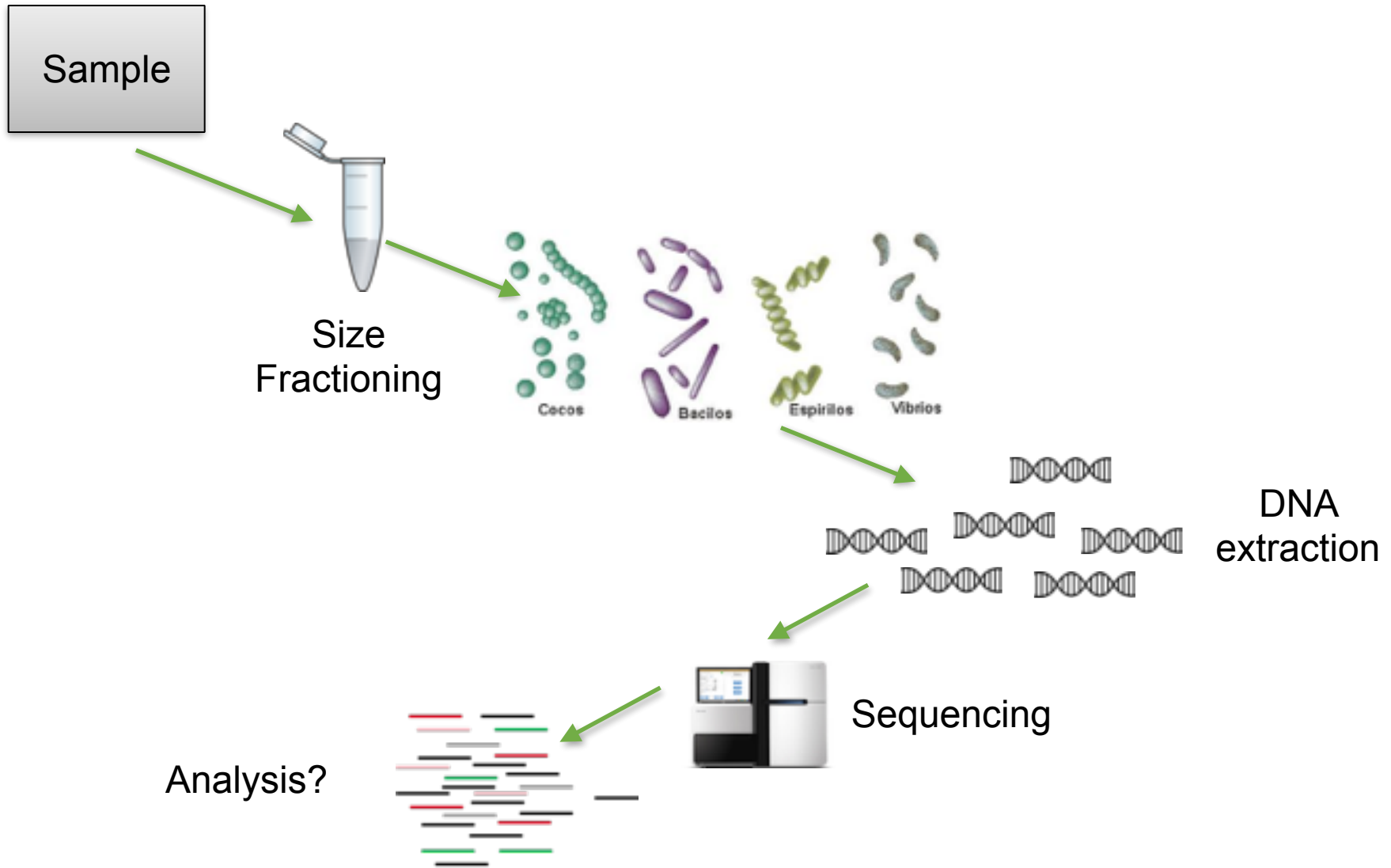
## Novel Gut-Based Pharmacology of Metformin in Patients with Type 2 Diabetes Mellitus

Antonella Napolitano<sup>1\*</sup>, Sam Miller<sup>2</sup>, Andrew W. Nicholls<sup>3</sup>, David Baker<sup>3</sup>, Stephanie Van Horn<sup>4</sup>, Elizabeth Thomas<sup>4</sup>, Deepak Rajpal<sup>5</sup>, Aaron Spivak<sup>5</sup>, James R. Brown<sup>5</sup>, Derek J. Nunez<sup>6</sup>

<sup>1</sup>Immuno-Inflammation Unit, GSK R&D, Stevenage, Herts, United Kingdom, <sup>2</sup>Quantitative Sciences, GSK R&D, Stevenage, Herts, United Kingdom, <sup>3</sup>Safety Assessment, GSK R&D, Ware, Herts, United Kingdom, <sup>4</sup>Target and Pathways Validation, GSK R&D, Upper Providence, Pennsylvania, United States of America, <sup>5</sup>Computational Biology, GSK R&D, Upper Providence, Pennsylvania, United States of America, <sup>6</sup>Endocrinology Discovery Unit, GlaxoSmithKline R&D, GSK R&D, Research Triangle Park, North Carolina, United States of America



# From environment to DNA sequence



# From environment to DNA sequence





**<http://www.ebi.ac.uk/metagenomics>**

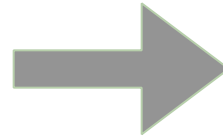


# Overview: EBI Metagenomics Portal

Easy submission



**Submission** of sequence data for archiving and analysis

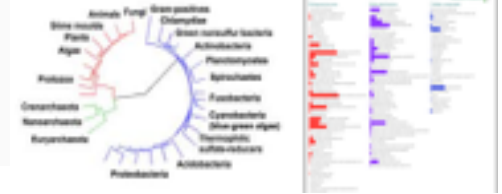


Powerful analysis



Data **analysis** using selected EBI and external software tools

Visualisation



Data presentation and **visualisation** through web interface



# Submitting to EBI Metagenomics

- EBI Metagenomics want to encourage people to supply as much detailed **metadata** as possible, but with the lowest possible overhead

who

where, when, what

how



name  
institute  
country ...  
contact



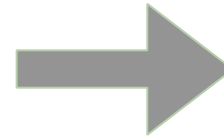
latitude and longitude,  
depth,  
salinity,  
temperature,  
time ...



quantity,  
conservation process,  
storage conditions,  
treatments,  
extraction methods ...



platform,  
protocol,  
filtering and QC,  
analysis,  
tool versions ...



EBI Metagenomics

- Development of intuitive web-based tools : **ENA Webin** and **ISA tools**
- Use of templates and check-lists (MIGS/MIXS standards)
- Tutorial and direct support





# Overview: EMG Portal analysis

Powerful analysis

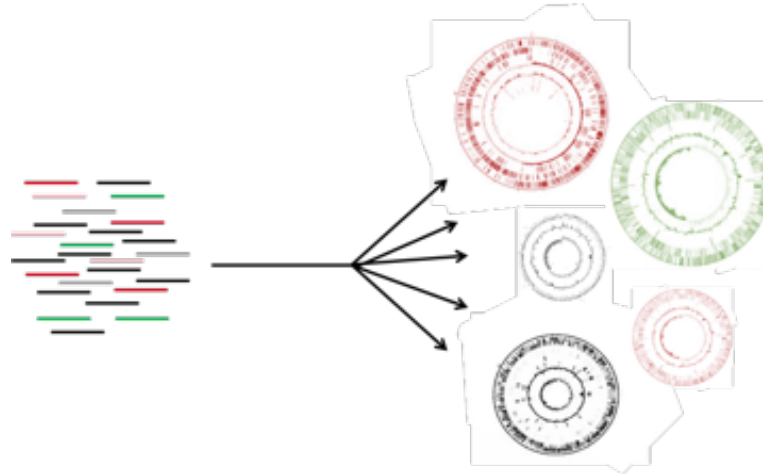


Data analysis using  
selected EBI and  
external software tools

- Provide robust sequence analysis services to all metagenomic researchers
  - Understand species diversity and functional potential of a community



# Metagenomics assembly?

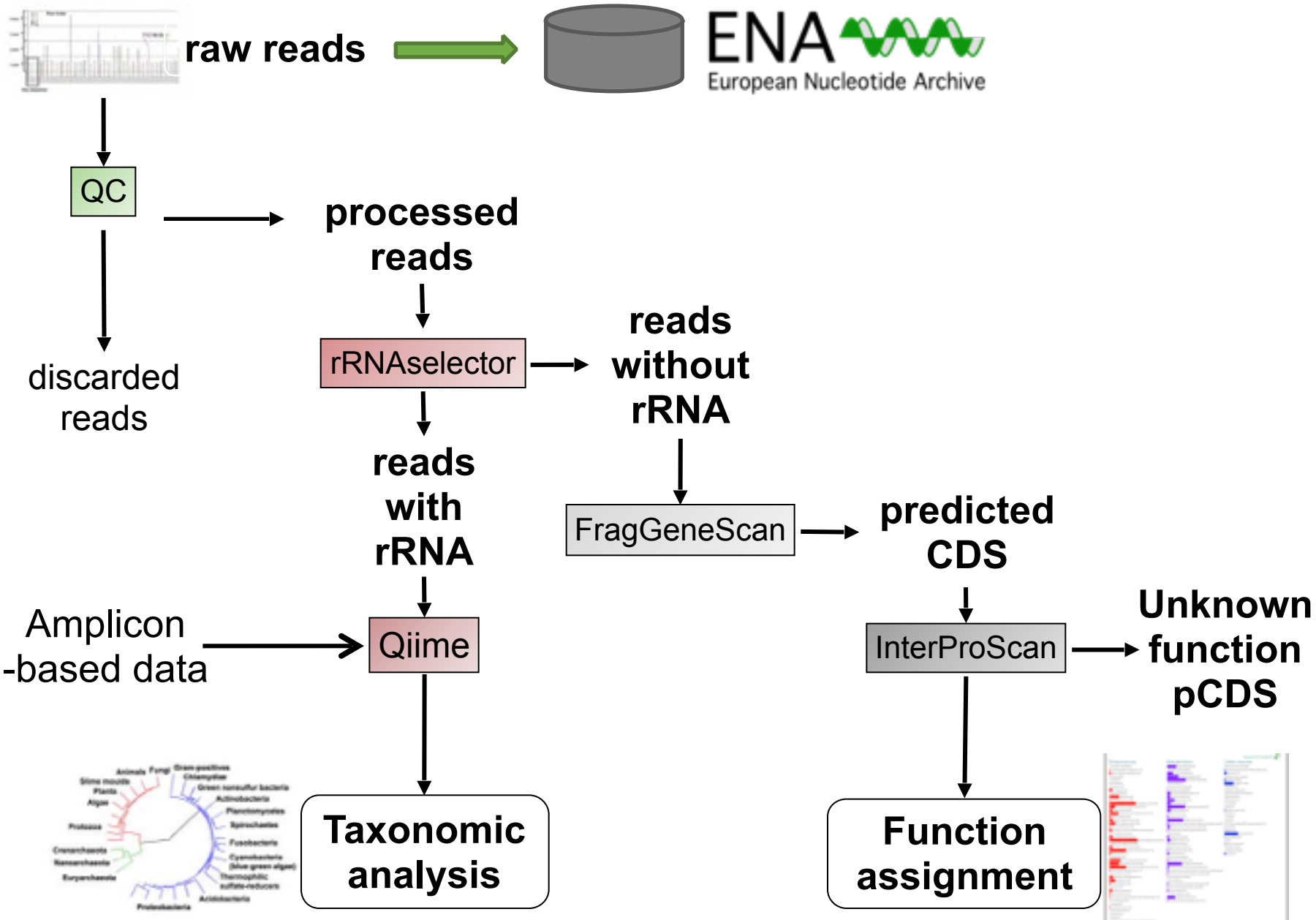


- Metagenomics: Not clear how you avoid assembling sequences from different species together : chimaera



- No reference sequence to align against





## EMG portal does not use blast based homology methods

Instead reads are compared to models (signatures) generated from multi-sequences alignments:

- more specific and meaningful annotations
- faster annotation
- rRNASelector identify 5, 16 and 28s rRNA (profile HMM models)  
=> 16s-based Qiime taxonomy annotations
- FragGenScan predict CDSs (HMM models)  
=> InterProScan functional annotations (profiles and models)



# Profile Hidden Markov Models

Input multiple alignment:

```
seq1  ACG-LD
seq2  SCG--E
Seq3  NCGgFD
Seq4  TCG-WQ
      123-45
```

Consensus columns assigned,  
Defining inserts and deletes:

# Profile Hidden Markov Models

Input multiple alignment:

```
seq1  ACG-LD
seq2  SCG--E
Seq3  NCGgFD
Seq4  TCG-WQ
      123-45
```

Consensus columns assigned,  
Defining inserts and deletes:

# Profile Hidden Markov Models

Input multiple alignment:

```
seq1  ACG-LD
seq2  SCG--E
Seq3  NCGgFD
Seq4  TCG-WQ
      123-45
```

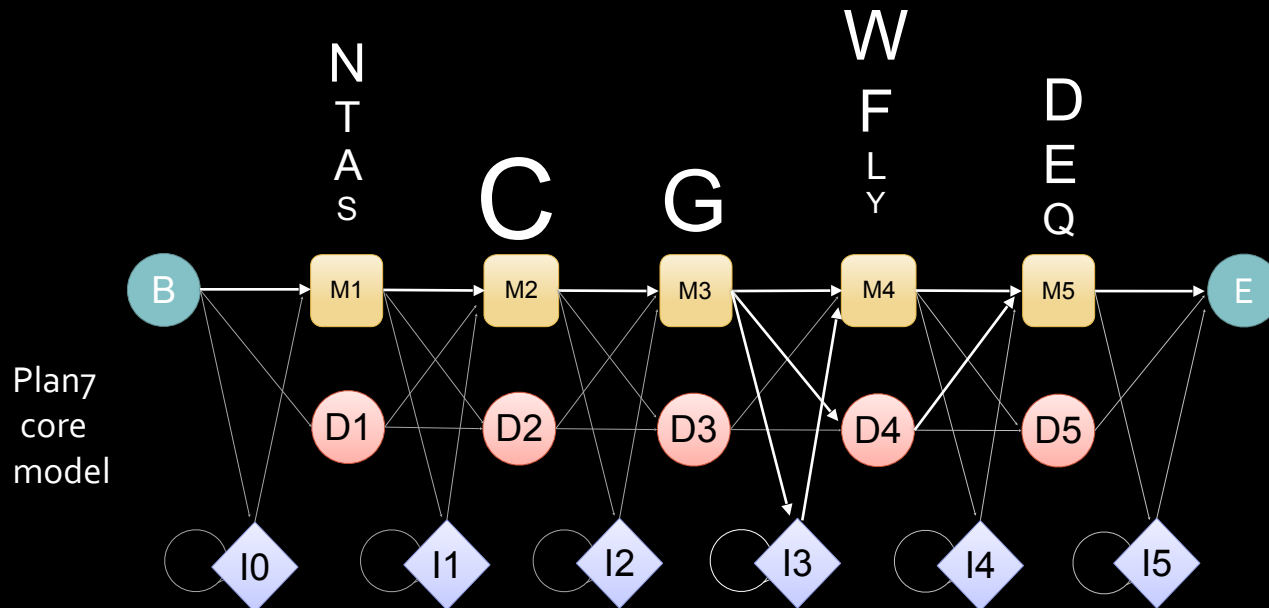
Consensus columns assigned,  
Defining inserts and deletes:

# Profile Hidden Markov Models

Input multiple alignment:

seq1 ACG-LD  
seq2 SCG--E  
Seq3 NCGgFD  
Seq4 TCG-WQ  
123-45

Consensus columns assigned,  
Defining inserts and deletes:



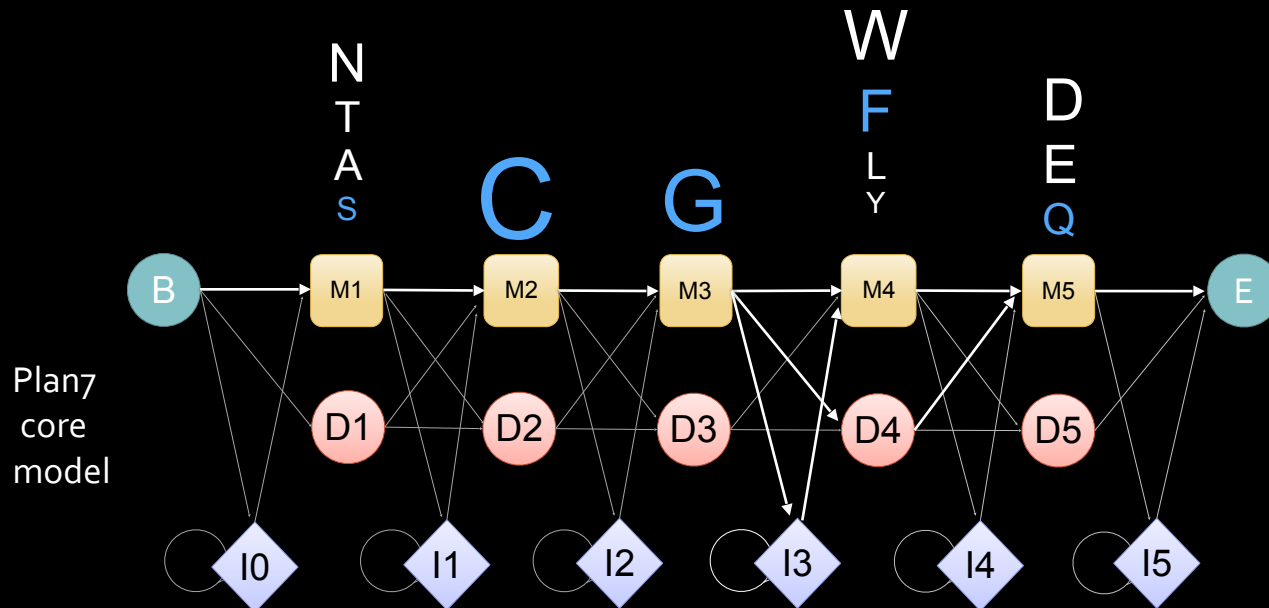


# Profile Hidden Markov Models

Input multiple alignment:

seq1 ACG-LD  
seq2 SCG--E  
Seq3 NCGgFD  
Seq4 TCG-WQ  
123-45

Consensus columns assigned,  
Defining inserts and deletes:



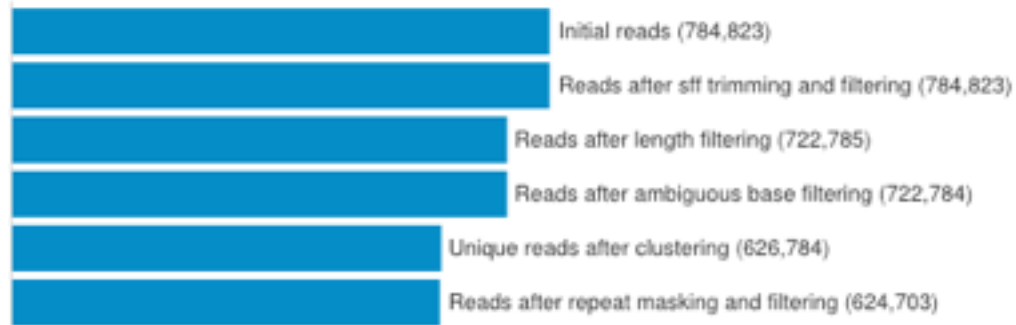
# QC processes



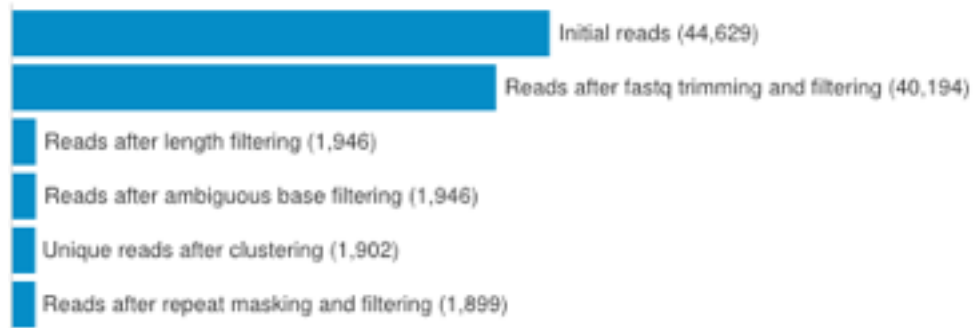
- **Clipping** - low quality ends trimmed and adapter sequences removed
- **Quality filtering** - sequences with  $> 10\%$  undetermined nucleotides removed
- **Read length filtering** - short sequences ( $< 100$  nt) are removed
- **Duplicate sequences removal** – clustered (99% identity UCLUST or 50 nt similarity Prefix) and representative sequence chosen
- **Repeat masking** - RepeatMasker (open-3.2.2), removes reads with 50% or more nucleotides masked (low complexity regions)

# QC effects by sequencing platform

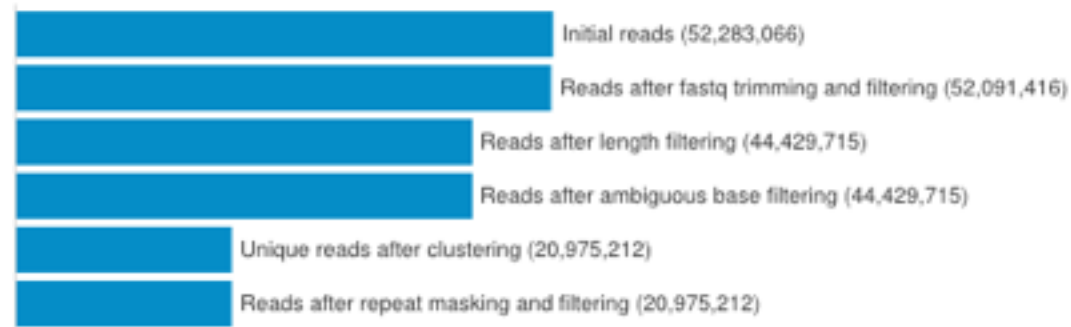
Roche 454



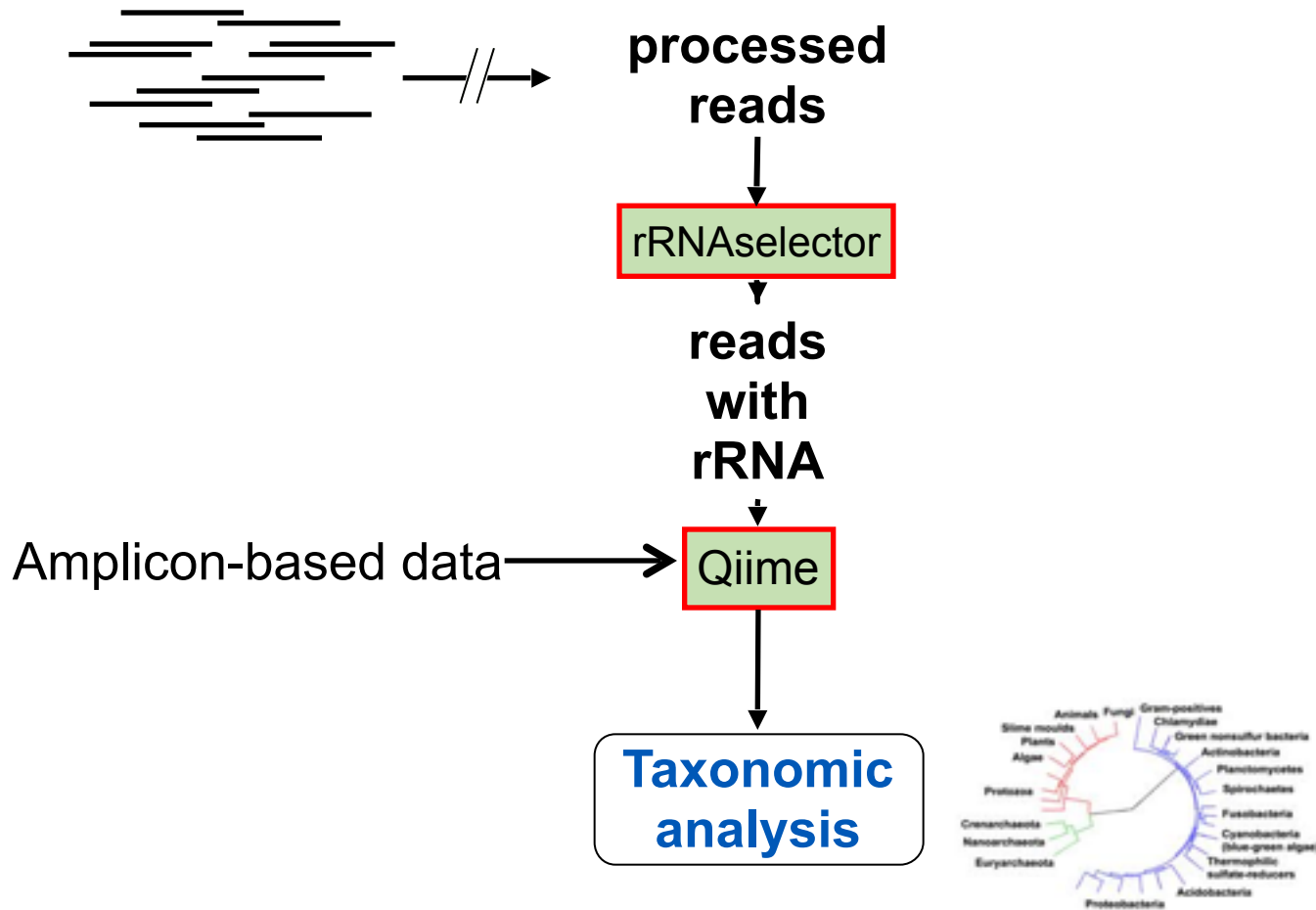
Ion Torrent



Illumina



# EBI Metagenomics: taxonomic analysis



# Common approaches to taxonomic analysis

- Identification of reads with 16S sequence (e.g. using rRNASelector) and closed-reference OTU picking in **QIIME**
- **Blast-based** analysis.
  - E.g. blasting reads against the NCBI non-redundant nucleotide or protein data databases and inferring taxonomic lineage from the best hit
  - The tool **MEGAN** requires Blast output. A major drawback is that without preprocessing of NGS datasets and access to a major computational resource, this is not an option for most.
- **MetaPhlAn** approach
  - (<http://huttenhower.sph.harvard.edu/metaphlan>)
  - relies on unique clade-specific marker genes identified from 3,000 reference genomes
  - fast, but limited to certain types of study (mainly human microbiome)



# Taxonomic analysis

Currently only taxonomy analysis for Prokaryotes

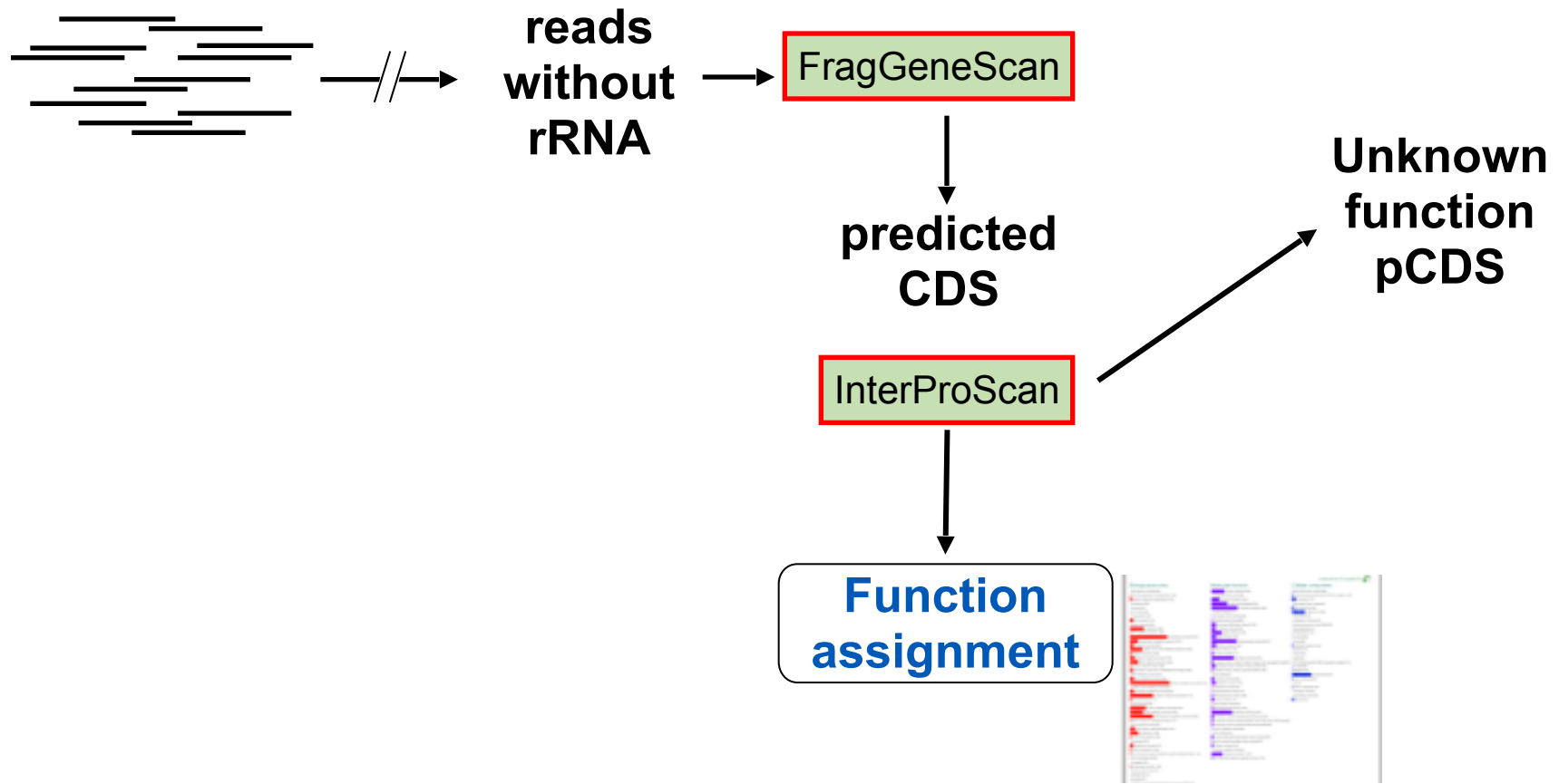
**rRNA sequences** are identified using [rRNASelector](#):

- hidden Markov models to identify rRNA sequences
- 60 bp minimum overlap with curated HMM model
- E-value <  $10^{-5}$

**Annotations** are associated using [Qiime](#):

- rRNA are annotated using the [Greengenes](#) reference database

# EBI Metagenomics: functional analysis



# EBI Metagenomics: functional annotation

The pipeline uses [FragGeneScan](#) to predict CDSs directly from the reads:

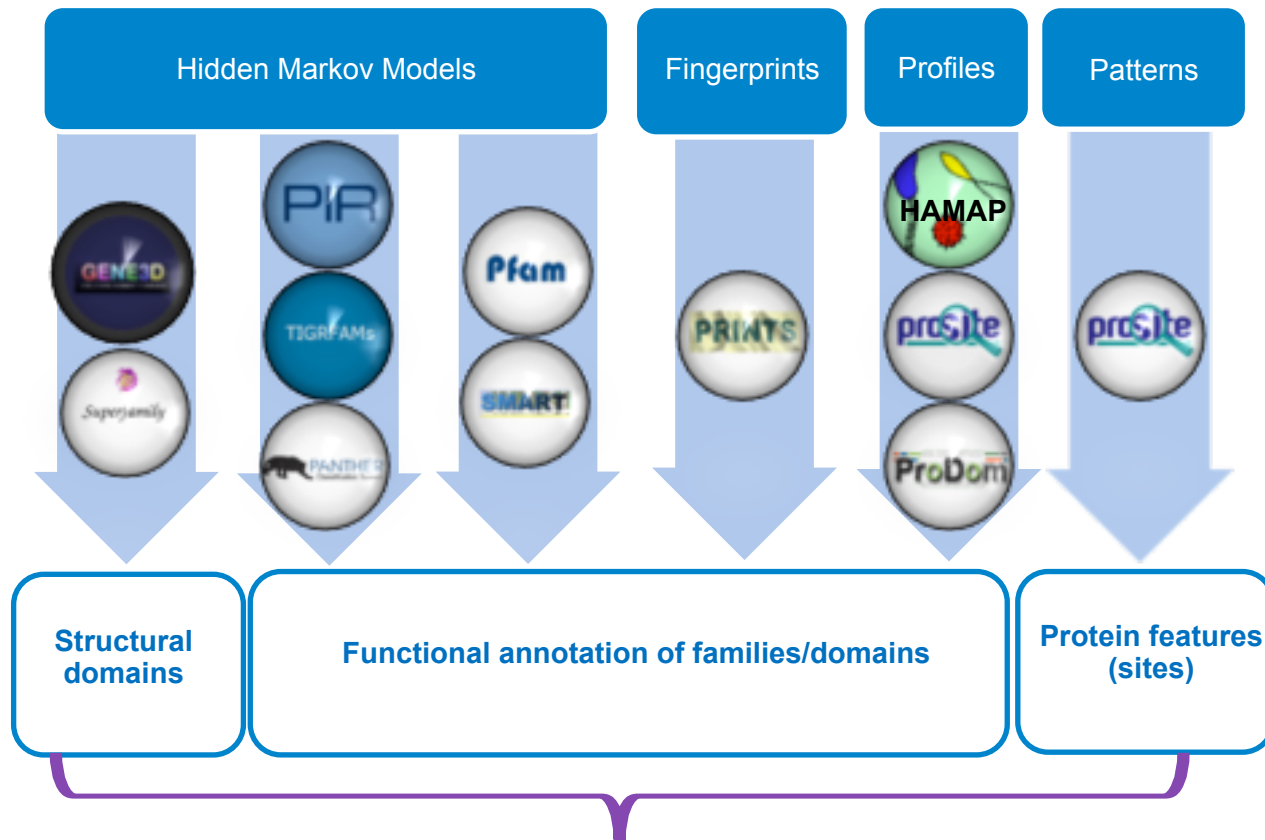
- hidden Markov models to correct frame-shift using codon usage
- probabilistic identification of start and stop codons
- 60 bp minimum ORF

Annotation is carried out using [InterProScan](#) with a subset of [InterPro's](#) databases

- analysis speed
- ability to cope with sequence fragments



# The benefits of InterPro



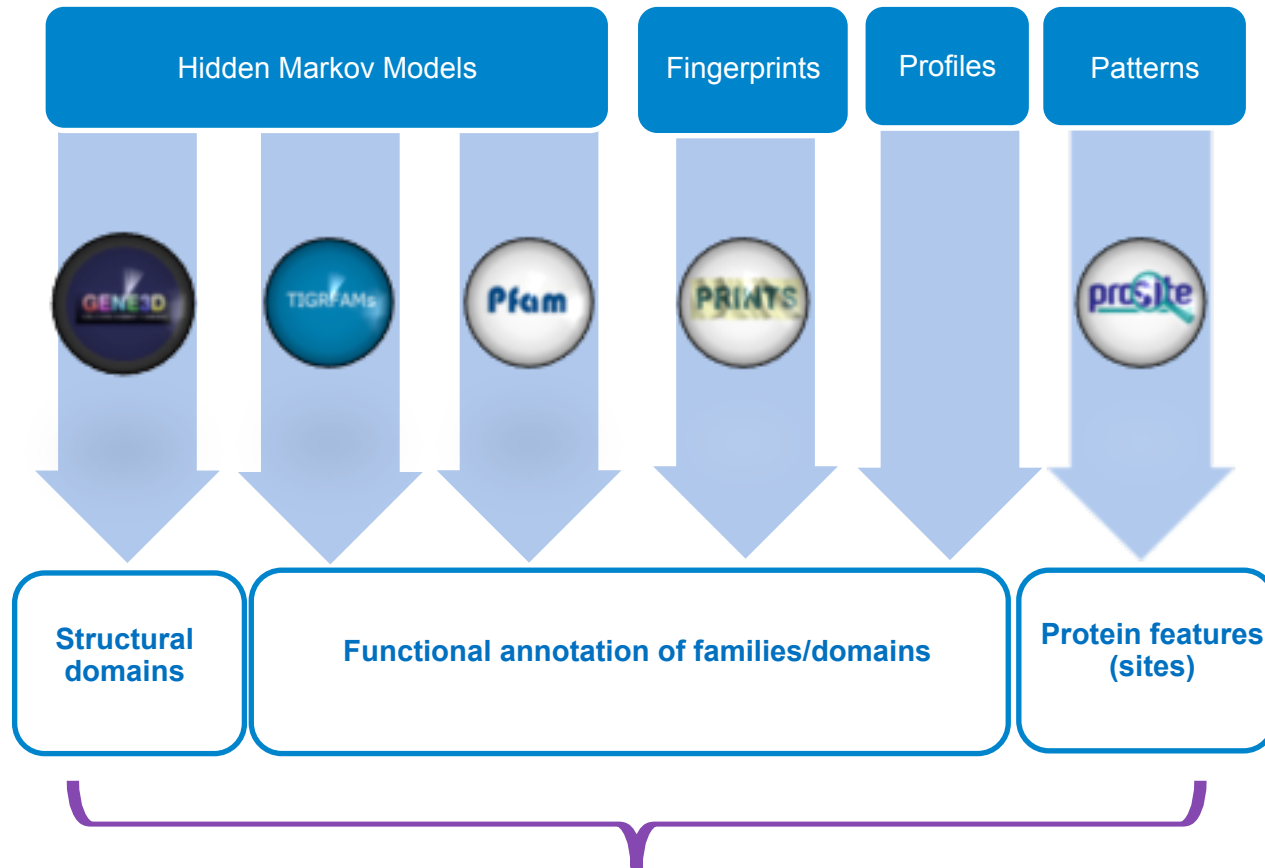
# Using InterPro for annotation

- Underlies automated systems that annotate UniProtKB/TrEMBL
- Provides matches to **90 million** proteins - **over 80%** of UniProtKB
- Source of **~ 170 million** GO mappings for **~ 50 million** distinct UniProtKB sequences

## Annotation consistency

- Using InterPro and GO allows **direct comparison** with proteins in UniProtKB

# InterPro in the Metagenomics Portal

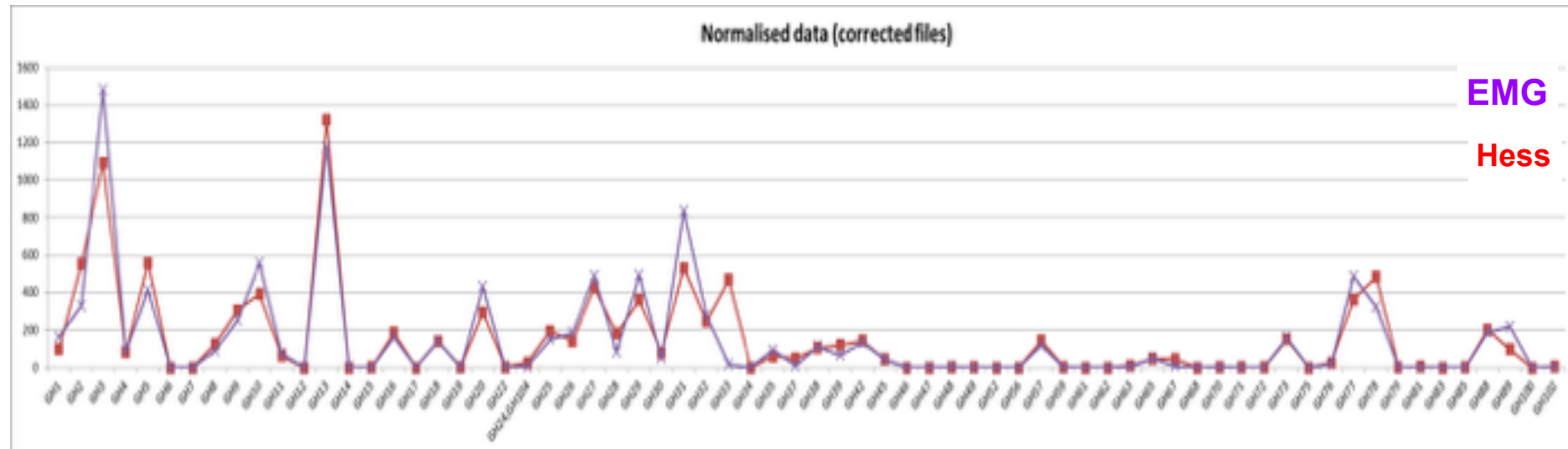


# Annotations without assembly

Re-analysis of Hess et al, Science (2011) 331:463

Metagenomic Discovery of Biomass-Degrading Genes and Genomes from Cow Rumen

Comparison of the normalised number of **genes** / **reads** corresponding to CAZy Glycoside Hydrolase Family from the **Hess et al** paper and from the **EMG pipeline**.



**Hess et al**: genome assembly then gene prediction using a subset of Pfam.

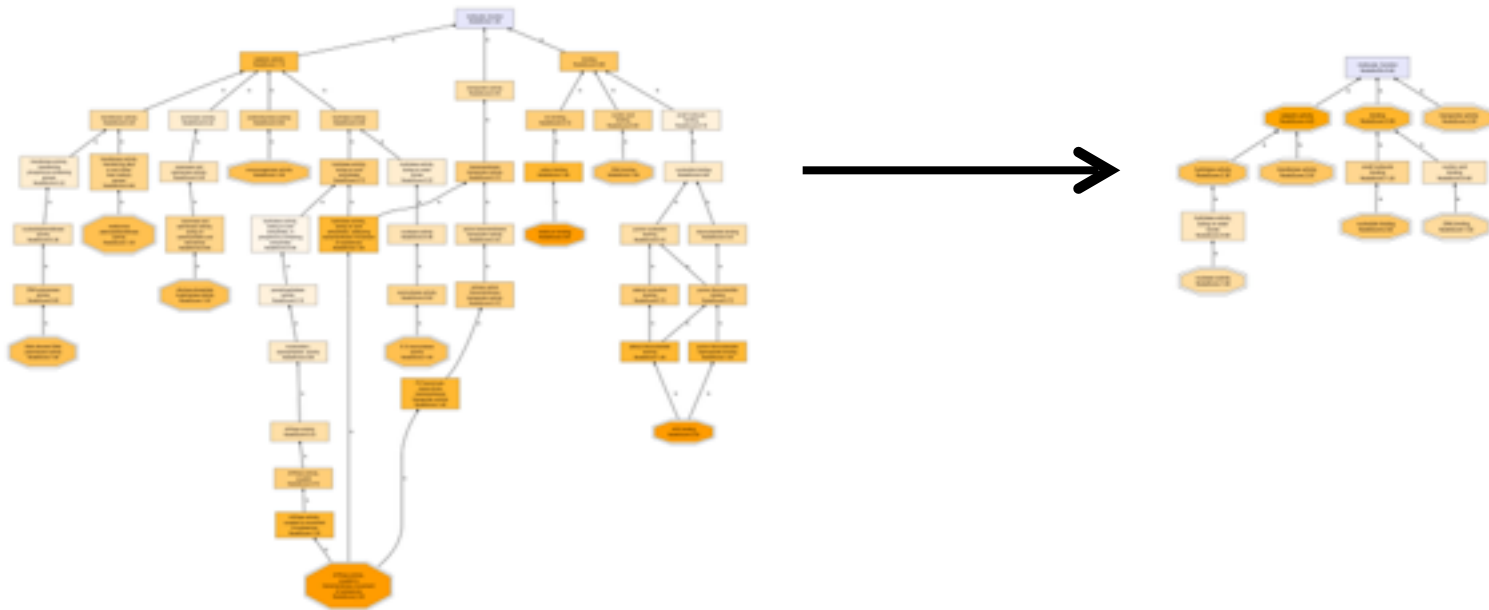
**EMG pipeline**: no assembly and gene prediction using InterPro.

Discrepancies are due to the different ways in which significance cut-off are calculated.

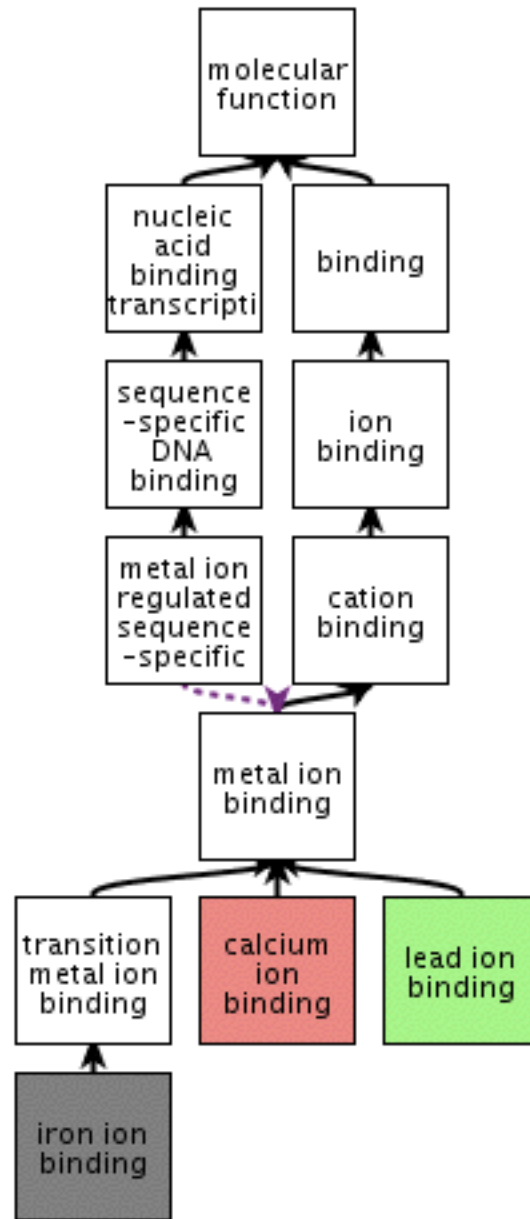


# Visualising data: GO Slims

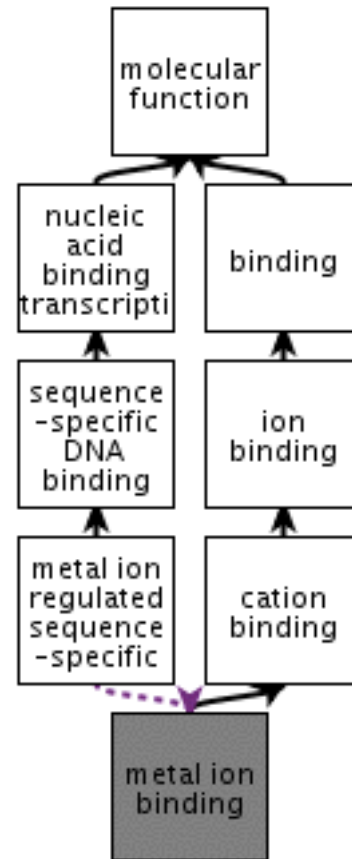
- Cut-down versions of the GO containing a subset of terms
- Give a broad overview of the ontology content without the detail of the specific fine-grained terms



# GO Slims



# GO Slims



Slimmed term:

# Metagenomics - Big Data

- Speed is really important! MAY
- Submitted nucleotide sequences: **27,509,856,436**
- Average length per sequence: **120 nt**
- Predicted CDS: **8,167,600,355**
- Total InterProScan matches: **1,866,871,818**
- Number of different samples: **2,808**
  
- **99.8%** of this data arrived and was processed in the last two years.  
**96.5%** of it is publicly available via the website.



# Metagenomics - Big Data

- Speed is really important! MONDAY
- Submitted nucleotide sequences: **43,315,534,332**
- Average length per sequence: **120 nt**
- Predicted CDS: **17,301,862,307**
- Total InterProScan matches: **3,276,195,744**
- Number of different samples: **4,330**
  
- **99.8%** of this data arrived and was processed in the last two years.  
**96.5%** of it is publicly available via the website.

# Downstream analysis: download options

Overview

Quality control

Taxonomy analysis

Functional analysis

Download

You can download in this section the full set of analysis results files and the original raw sequence reads.

## Sequence data

- Submitted nucleotide reads (ENA website)
- Processed nucleotide reads (FASTA) - 2 MB
- Processed reads with pCDS (FASTA) - 2 MB
- Processed reads with InterPro matches (FASTA) - 1 MB
- Processed reads without InterPro match (FASTA) - 835 KB
- Predicted CDS (FASTA) - 710 KB
- Predicted CDS with InterPro matches (FASTA) - 451 KB

## Functional Analysis

- InterPro matches (TSV) - 1 MB
- Complete GO annotation (CSV) - 44 KB
- GO slim annotation (CSV) - 7 KB

## Taxonomic Analysis

- Reads encoding 5S rRNA (FASTA) - 565 bytes
- Reads encoding 16S rRNA (FASTA) - 21 KB
- Reads encoding 23S rRNA (FASTA) - 37 KB
- OTUs and taxonomic assignments (BIOM) ⓘ - 6 KB
- Phylogenetic tree (Newick format) ⓘ - 289 bytes
- OTUs and taxonomic assignments (TSV) - 2 KB

**relatively small result files: can be used for downstream analysis with other tools**

# Overview: EMG Portal output

## Visualisation



Data presentation and visualisation through web interface

- Assist laboratory researchers handle and make sense of massive volumes of sequence data
  - Do this by designing intuitive, user-friendly web interfaces
  - Browse, visualise and download





EBI Metagenomics &gt; Projects

## Projects list

Text:

Privacy:

Public projects can be browsed and searched using names and keywords

1 - 10 of 41


[Download detailed info \(CSV\)](#) [Download table \(CSV\)](#)

Project name	Samples	Submitted date
A core gut microbiome in obese and lean twins	18	03-Mar-2010
A metagenomics transect into the deepest point of the Baltic Sea reveals clear stratification of microbial functional capacities	0	27-Sep-2013
Antarctica Aquatic Microbial Metagenome	18	12-Jul-2011
Arctic Winter marine ecosystem	1	19-Mar-2013
BGI Type 2 Diabetes study (SRP008047)	145	14-Mar-2013
BGI Type 2 Diabetes study (SRP011011)	218	26-Mar-2013
Beta Lactam Antibiotics and Human Gut Microbiota	12	23-Sep-2013
Brazos-Trinity Basin Sediment Metagenome: IODP Site 1320	1	04-Dec-2013
Buffalo rumen metagenomics	1	20-Mar-2013
Comparative freshwater metagenomics of Swedish and American lakes	0	04-Dec-2013





EBI Metagenomics - Project: A core gut microbiome in obese and lean twins

## Project overview (SRP000319)

## A core gut microbiome in obese and lean twins

BioProject ID: #32089

Submitted date: 03-Mar-2010

## Description

We have characterized the fecal microbial communities of adult female monozygotic and dizygotic twin pairs concordant for leanness or obesity, and their mothers. The results demonstrate that a diversity of organismal assemblages can nonetheless yield a core microbiome at a functional level, and that deviations from this core are associated with different physiologic states (obese versus lean)

Experimental factor: obese vs lean

## Contact details

Institute: COS-GL

Name: Turnbaugh PJ

## Associated samples

Sample name	Sample ID	Collection date	Source	Analysis results
Obese human gut (patient T529)	SRS009826	-	Host associated	<a href="#">Taxonomy</a>   <a href="#">Function</a>   <a href="#">J.S.</a>
Obese human gut (patient T528)	SRS009826	-	Host associated	<a href="#">Taxonomy</a>   <a href="#">Function</a>   <a href="#">J.S.</a>
Lean human gut (patient T53)	SRS000999	-	Host associated	<a href="#">Taxonomy</a>   <a href="#">Function</a>   <a href="#">J.S.</a>
Lean human gut (patient T51)	SRS000998	-	Host associated	<a href="#">Taxonomy</a>   <a href="#">Function</a>   <a href="#">J.S.</a>
Obese human gut (patient T519)	SRS001007	-	Host associated	<a href="#">Taxonomy</a>   <a href="#">Function</a>   <a href="#">J.S.</a>
Overweight human gut (patient T53)	SRS001006	-	Host associated	<a href="#">Taxonomy</a>   <a href="#">Function</a>   <a href="#">J.S.</a>
Lean human gut (patient T58)	SRS001005	-	Host associated	<a href="#">Taxonomy</a>   <a href="#">Function</a>   <a href="#">J.S.</a>
Lean human gut (patient T57)	SRS001004	-	Host associated	<a href="#">Taxonomy</a>   <a href="#">Function</a>   <a href="#">J.S.</a>
Obese human gut (patient T56)	SRS001003	-	Host associated	<a href="#">Taxonomy</a>   <a href="#">Function</a>   <a href="#">J.S.</a>
Lean human gut (patient T55)	SRS001002	-	Host associated	<a href="#">Taxonomy</a>   <a href="#">Function</a>   <a href="#">J.S.</a>

Overall project description

Project-associated publication(s)

## Related Publications


- A core gut microbiome in obese and lean twins. Turnbaugh PJ, Hamady M, Wotjenko T. 2009 Nature. 2009 Jan 22;457(7226):480-4. Epub 2008 Nov 30.

Samples in the study

Direct link to analysis results



EMBL-EBI Services Research Training About us

 EBI Metagenomics

Home Submit data Projects **Samples** About EBI Metagenomics Contact

EBI Metagenomics - Project: A core gut microbiome in obese and lean twins - Sample: Obese human gut (patient TS28)

**Sample** (SRP009825)  
**Obese human gut (patient TS28)**

Overview **Quality control** Taxonomy analysis Functional analysis Download

**Description**

Sample Characteristics - The MOAFTS twin cohort, comprised of female like-sex twin pairs, were identified from Missouri birth records over the period 1994-1999, when the twins had a median age of 15 years. A total of 350 twins from the larger MOAFTS cohort completed screening interviews for the present study. We were able to take advantage of the wave five interview of the MOAFTS twin cohort (which has 90% retention of wave four participants) to identify pairs most likely to meet study criteria. Eligibility was then confirmed at screening interview.

Classification: Host-associated > Human > Digestive system > Large intestine > Fecal  
 Project: [A core gut microbiome in obese and lean twins \(SRP000319\)](#)

**Host associated**

Species: [Homo sapiens](#) #Tax ID 9606  
 Sex: Female  
 Phenotype: Obese

**Localisation**

Geographic location: United States of America: Missouri

**Other information**

Sex	female
Geographic location (longitude)	not available
Geographic location (country and/or sea_region)	USA; Missouri
Collection date	not available
Environment (biome)	human microbiome
Environment (feature)	human-associated habitat ENVO:00009003
Environment (material)	faeces ENVO:00002003
Sequencing method	pyrosequencing
Family relationship	mono-zygotic twin
Host common name	Homo sapiens
Host subject id	F10T10e1
Host taxid	9606
Miscellaneous parameter	period without antibiotics; 16 months
Phenotype	obese
Geographic location (latitude)	not available
Alt id - submitted sample id	TS28
GOLD sample classification	Host-associated > Human > Digestive system > Large intestine > Fecal
NCBI sample classification	77133

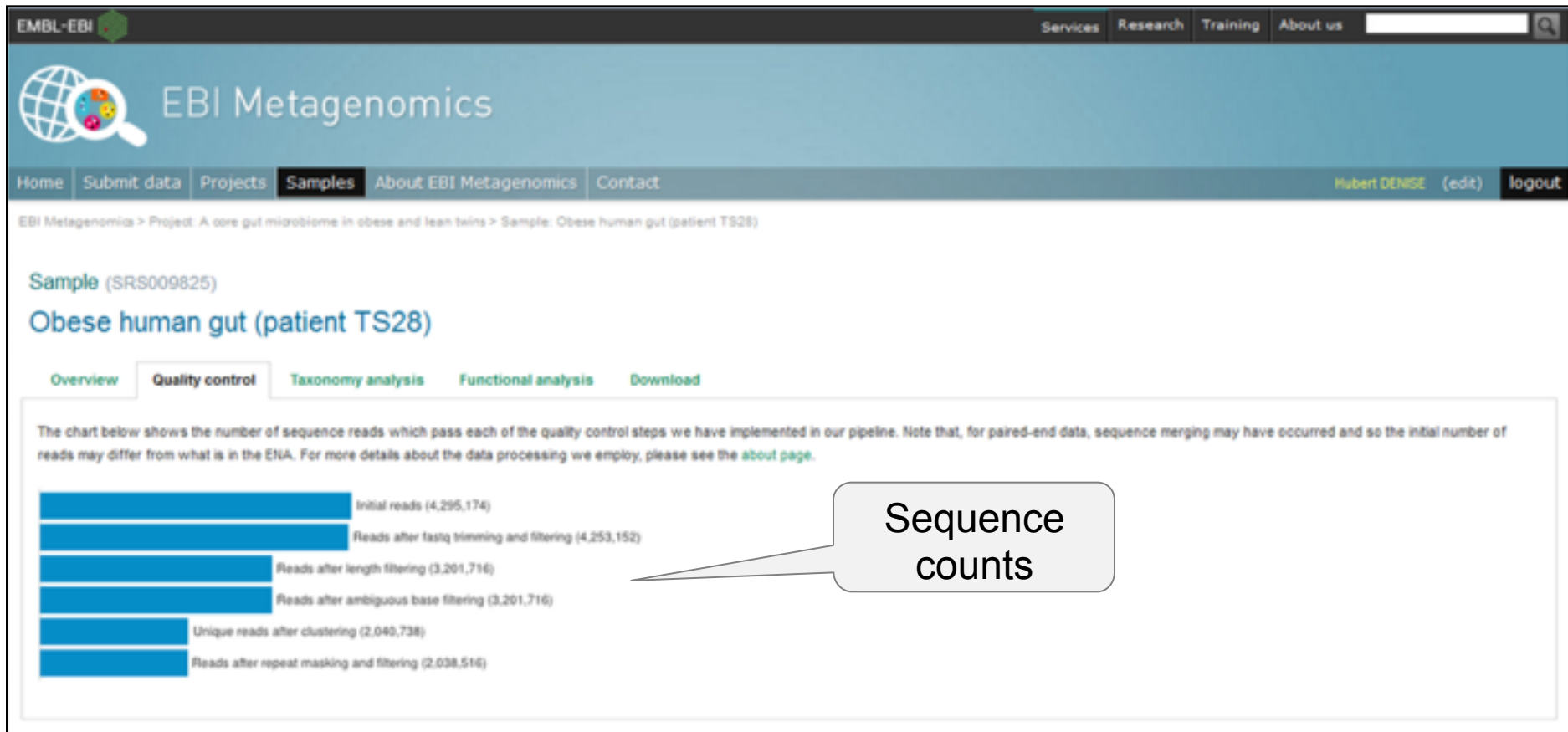
QC, analysis results and download tabs

Data in related resources

Descriptive meta-data



# EBI Metagenomics: QC tab



The screenshot shows the EBI Metagenomics website interface. At the top, there is a navigation bar with links for Services, Research, Training, and About us. Below this is the EBI Metagenomics logo and a search bar. A secondary navigation bar includes Home, Submit data, Projects, Samples (highlighted), About EBI Metagenomics, and Contact. The user is logged in as Hubert DENISE (edit) and has a logout button.

The main content area displays the sample information: **Sample (SRS009825)** **Obese human gut (patient TS28)**. Below the sample name are tabs for Overview, Quality control (selected), Taxonomy analysis, Functional analysis, and Download.

A text block explains the chart: "The chart below shows the number of sequence reads which pass each of the quality control steps we have implemented in our pipeline. Note that, for paired-end data, sequence merging may have occurred and so the initial number of reads may differ from what is in the ENA. For more details about the data processing we employ, please see the [about page](#)."

The chart is a horizontal bar chart showing the following sequence counts:

Quality Control Step	Number of Reads
Initial reads	4,295,174
Reads after fastq binning and filtering	4,253,152
Reads after length filtering	3,201,716
Reads after ambiguous base filtering	3,201,716
Unique reads after clustering	2,040,738
Reads after repeat masking and filtering	2,038,516

A callout box labeled "Sequence counts" points to the chart.



# EBI Metagenomics: taxonomy analysis tab

Google charts dynamic representation

Switch to bar chart, column or Krona interactive views

Export charts



Overview Quality control Taxonomy analysis Functional analysis Download

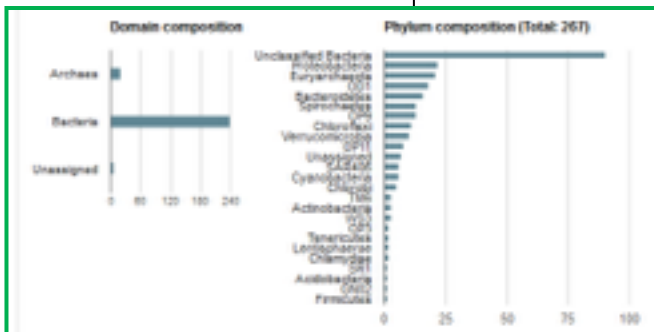
Top taxonomy Hits

Switch view: Export

**Domain composition**

**Phylum composition (Total: 267)**

Filter:	Domain	Unique OTUs	% unique OTUs	Count of reads assigned	% reads assigned
1	Unclassified Bacteria	90	33.71	0	0
2	Proteobacteria	22	8.24	0	0
3	Euryarchaeota	21	7.67	0	0
4	OD1	18	6.74	0	0
5	Bacteroidetes	16	5.99	0	0
6	Spirochaetes	13	4.87	0	0
7	OP9	13	4.87	0	0
8	Chloroflexi	11	4.12	0	0
9	Verrucomicrobia	10	3.75	0	0
10	OP11	8	3	0	0





# EBI Metagenomics: functional analysis tab

Google charts dynamic representation

Links to InterPro website

Export charts

Switch to bar chart, view

The screenshot displays the 'Functional analysis' tab of the EBI Metagenomics interface. It features a navigation bar with 'Overview', 'Quality control', 'Taxonomy analysis', 'Functional analysis', and 'Download'. The main content is divided into three sections:

- InterPro match summary:** Shows 'Most frequently found InterPro matches to this sample:' with an 'Export' button. Below is a pie chart and a table of matches.
- GO Terms annotation:** States 'A summary of Gene Ontology (GO) terms derived from InterPro matches to your sample is provided in the charts below.' It includes a 'Switch view:' section with a pie chart icon selected and a bar chart icon, along with an 'Export' button.
- Biological process, Molecular function, and Cellular component:** Each section contains a pie chart and a legend of terms.

Rank	Match	ID	Hits
1	NAD(P)-binding domain	IPRO10040	5540
2	Outer membrane insertion C-terminal signal, omp85 target	IPRO17090	3572
3	Aldolase-type TIM barrel	IPRO13785	3400
4	Rosemann-like alpha/beta/alpha sandwich fold	IPRO14729	2674
5	Pyridoxal phosphate-dependent transferase, major region, subdo	IPRO15421	2075
6	Tetrahiopeptide-like helical	IPRO11990	1855
7	Winged helix-turn-helix transcription repressor DNA-binding	IPRO11991	1631
8	ATPase-like, ATP-binding domain	IPRO03594	1493
9	ABC transporter-like	IPRO03439	1484
10	Signal transduction response regulator, receiver domain	IPRO01789	1354

**InterPro matches summary (Total: 5452)**

- NAD(P)-binding domain
- Outer membrane insertion C-terminal signal, omp85 target
- Aldolase-type TIM barrel
- Rosemann-like alpha/beta/alpha sandwich fold
- Pyridoxal phosphate-dependent transferase, major region, subdo
- Tetrahiopeptide-like helical
- Winged helix-turn-helix transcription repressor DNA-binding
- ATPase-like, ATP-binding domain
- ABC transporter-like
- Other

**Biological process**

- nitrogen compound metabolic proc
- biosynthetic process
- oxidation-reduction process
- small molecule metabolic process
- protein metabolic process
- transport
- RNA metabolic process
- cellular amino acid metabolic proce
- DNA metabolic process
- Other

**Molecular function**

- nucleotide binding
- oxidoreductase activity
- hydrolase activity
- transferase activity
- nucleic acid binding
- cofactor binding
- transporter activity
- ligase activity
- metal ion binding
- Other

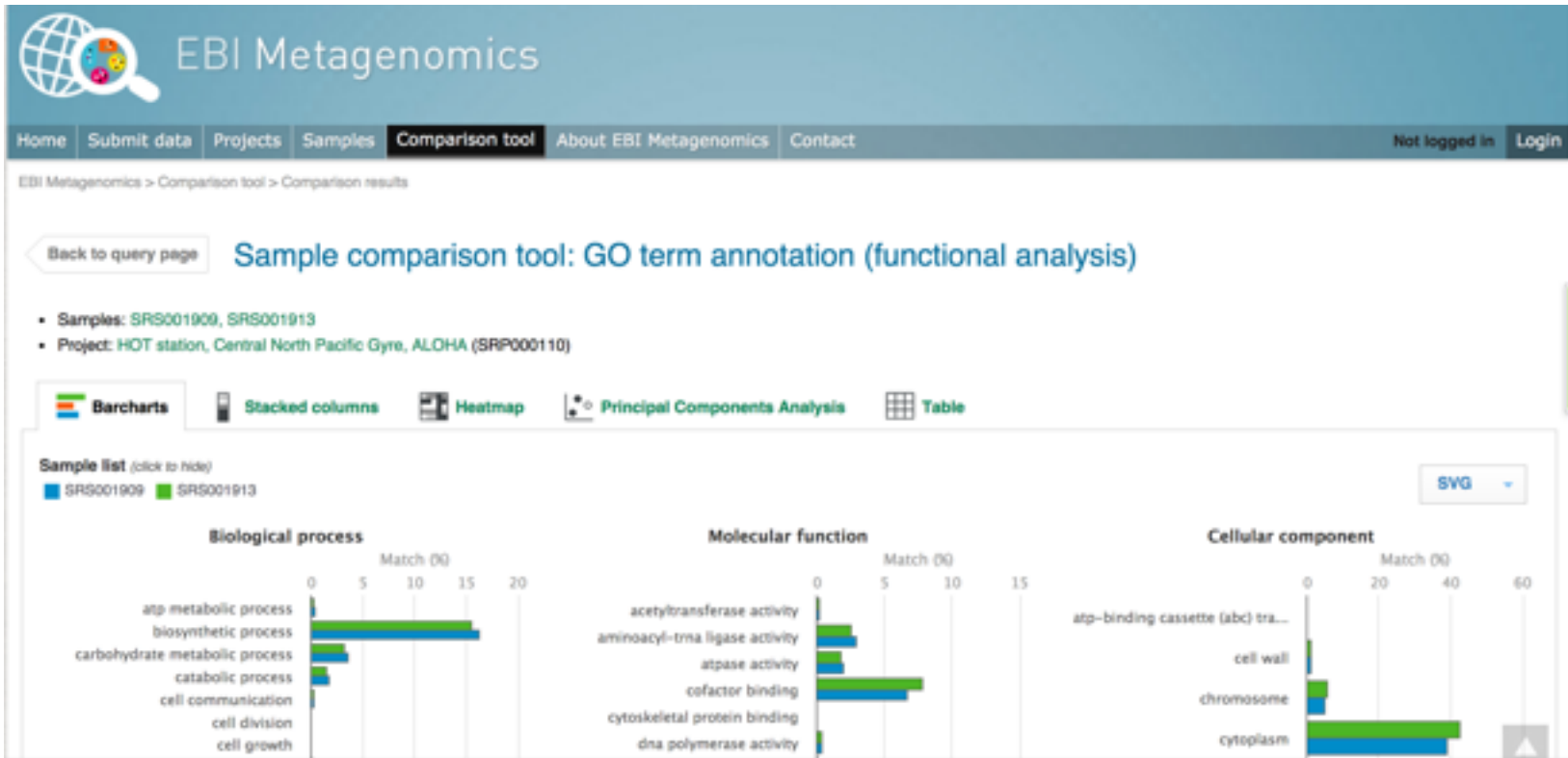
**Cellular component**

- membrane
- cytoplasm
- ribosome
- chromosome
- periplasmic space
- flagellum
- outer membrane
- plasma membrane
- Other



# Sample Comparisons

- <https://www.ebi.ac.uk/metagenomics/compare>



# Sample

- <https://www.ebi.ac.uk/metagenomics/compare>

EBI Metagenomics

Home Submit data Projects

EBI Metagenomics > Comparison tool > C

Back to query page **Sam**

- Samples: SRS001909, SRS001913
- Project: HOT station, Central N

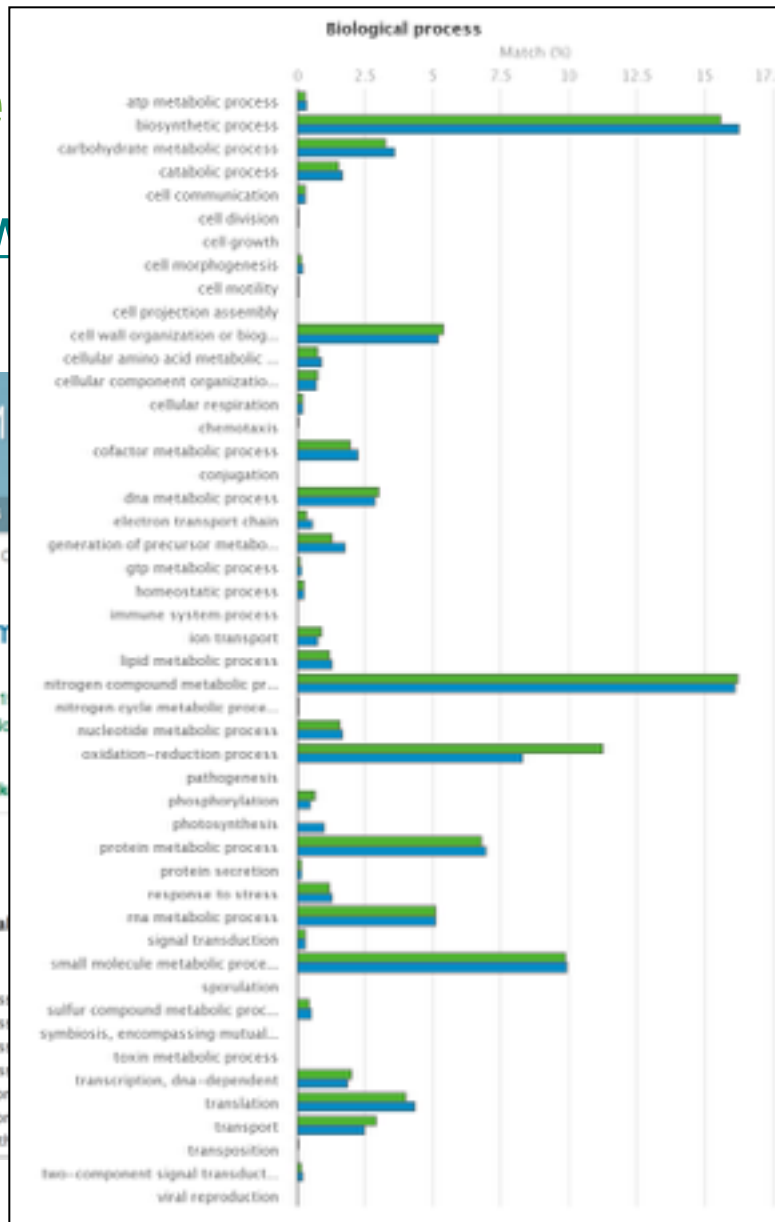
Barcharts Stack

Sample list (click to hide)

■ SRS001909 ■ SRS001913

Biological

- atp metabolic process
- biosynthetic process
- carbohydrate metabolic process
- catabolic process
- cell communication
- cell division
- cell growth



omics/compare

Not logged in Login

analysis)

Cellular component

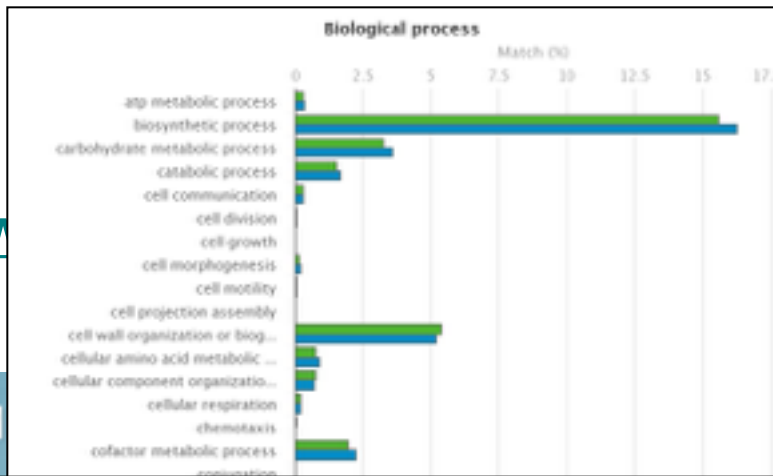
Match (%)

Cellular component	SRS001909 (Match %)	SRS001913 (Match %)
atp-binding cassette (abc) tra...	0.5	0.5
cell wall	0.5	0.5
chromosome	1.5	1.5
cytoplasm	40.5	40.5



# Sample

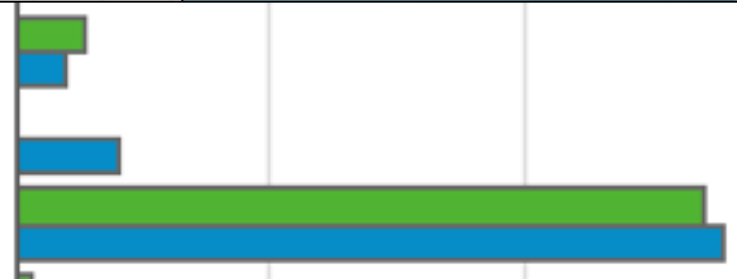
- <https://www.ebi.ac.uk/ena/browser/view/SRS001909>



<https://www.ebi.ac.uk/ena/browser/view/SRS001909>



phosphorylation  
photosynthesis  
protein metabolic process



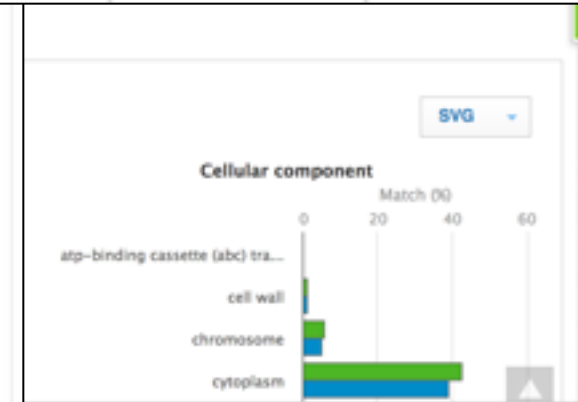
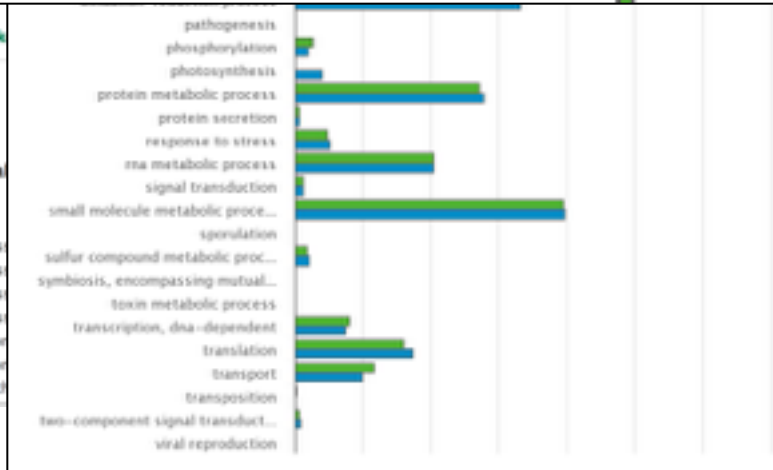
**Barcharts** **Stack**

Sample list (click to hide)

■ SRS001909 ■ SRS001913

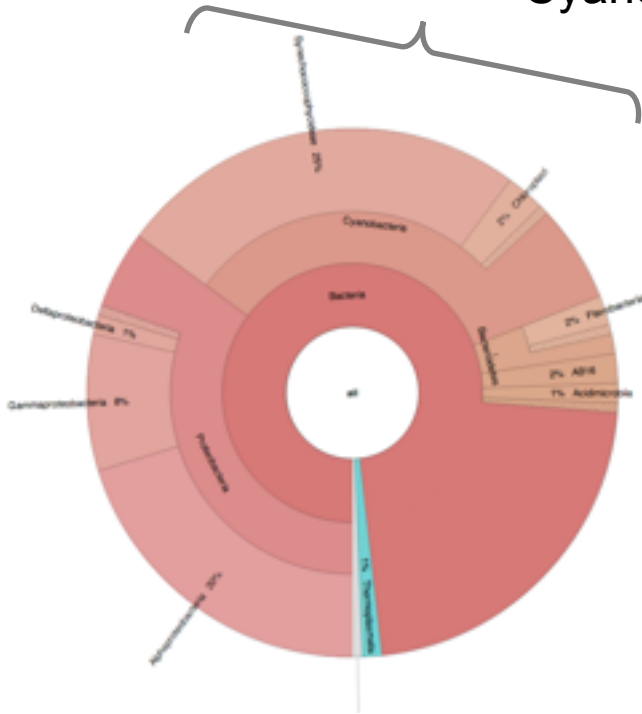
**Biological**

- atp metabolic process
- biosynthetic process
- carbohydrate metabolic process
- catabolic process
- cell communication
- cell division
- cell growth

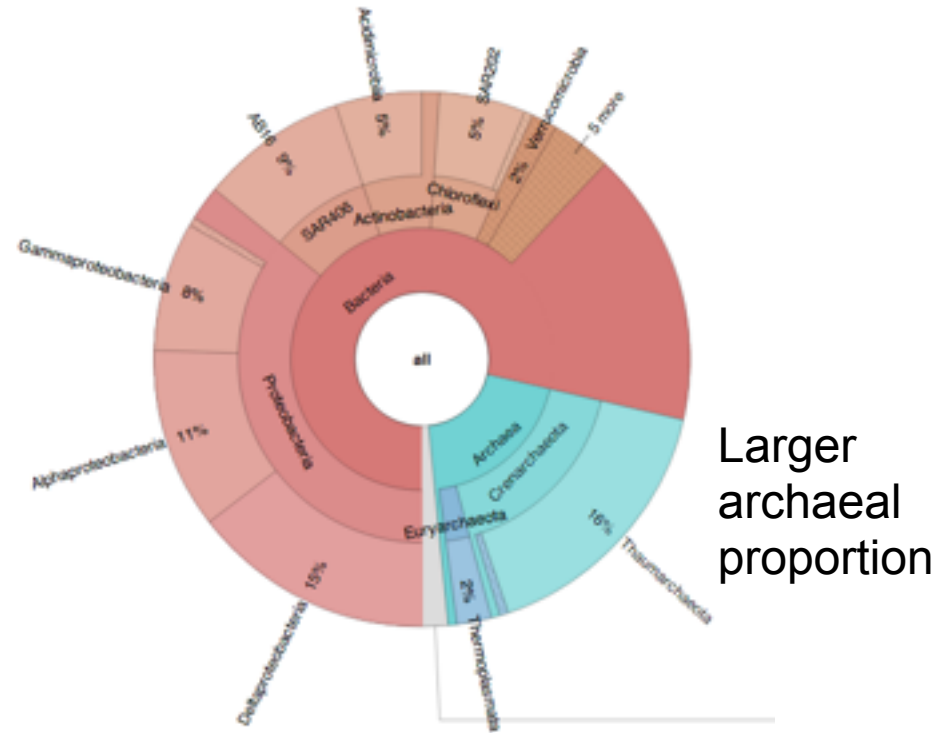


# Comparison of two Marine Biomes - Taxonomic distributions

Cyanobacteria



25m Depth



500m Depth

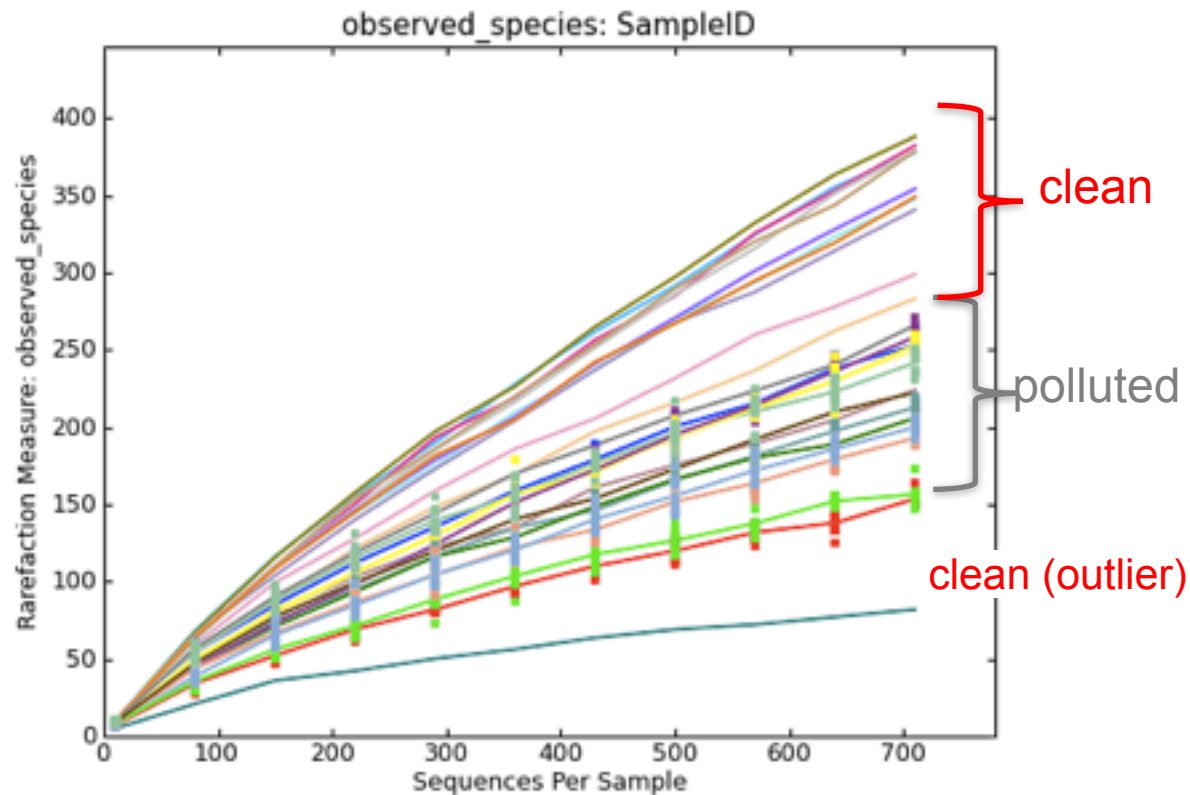


# EBI Metagenomics: application of taxonomy analysis

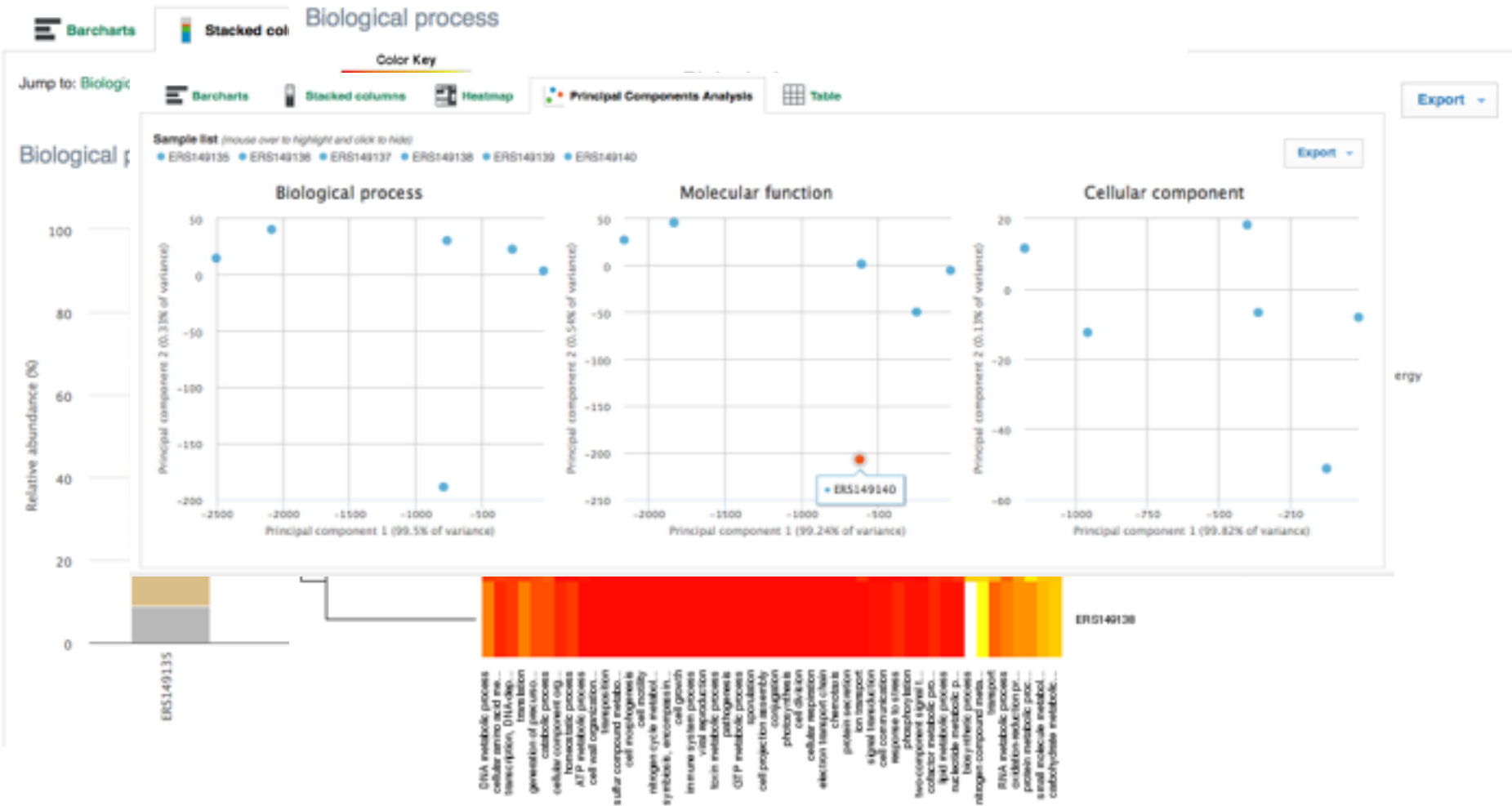
Sutton et al, Appl. Environ. Microbiol (2013), 79(2):619

Impact of Long-Term Diesel Contamination on Soil Microbial Community Structure.

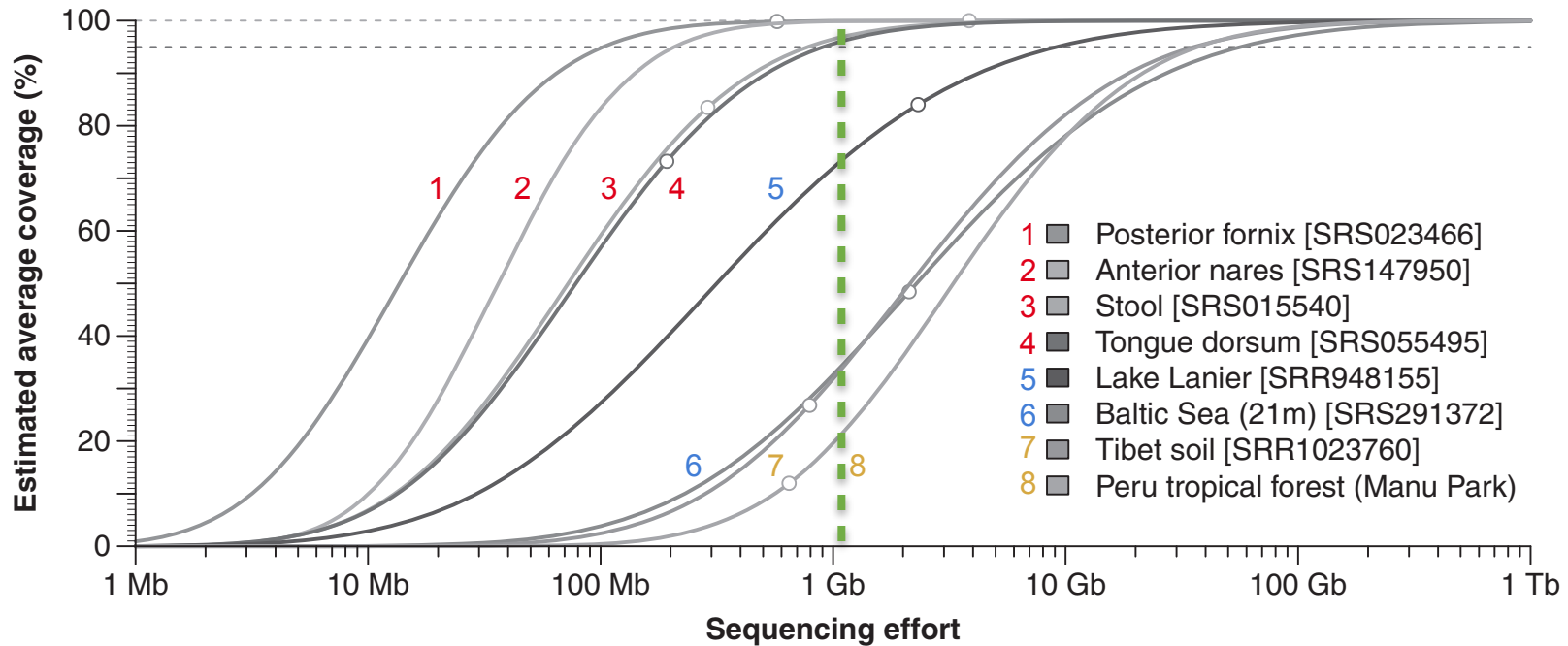
Alpha diversity analysis



# Comparison tool



# Show much of the microbial community has been sequenced?

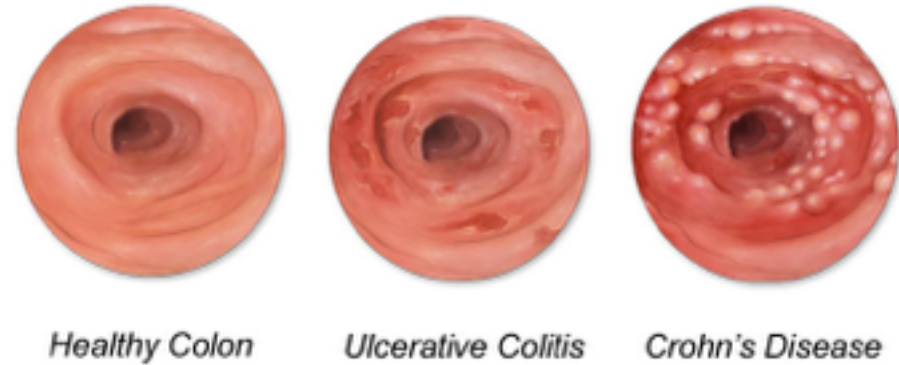
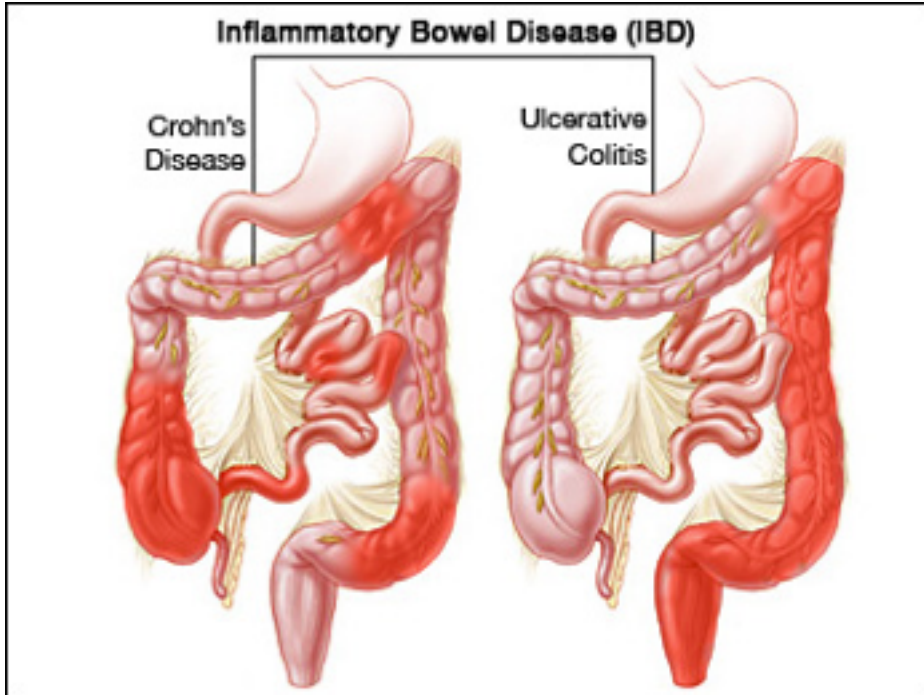


Rodriguez-R and Konstantinidis, *ISMEJ*, 2014 1-3





# Metagenomics - Human gut microbiome



Susceptibility to inflammatory bowel disease (IBD) :  
Friedreich ataxia (frataxin; fxn)

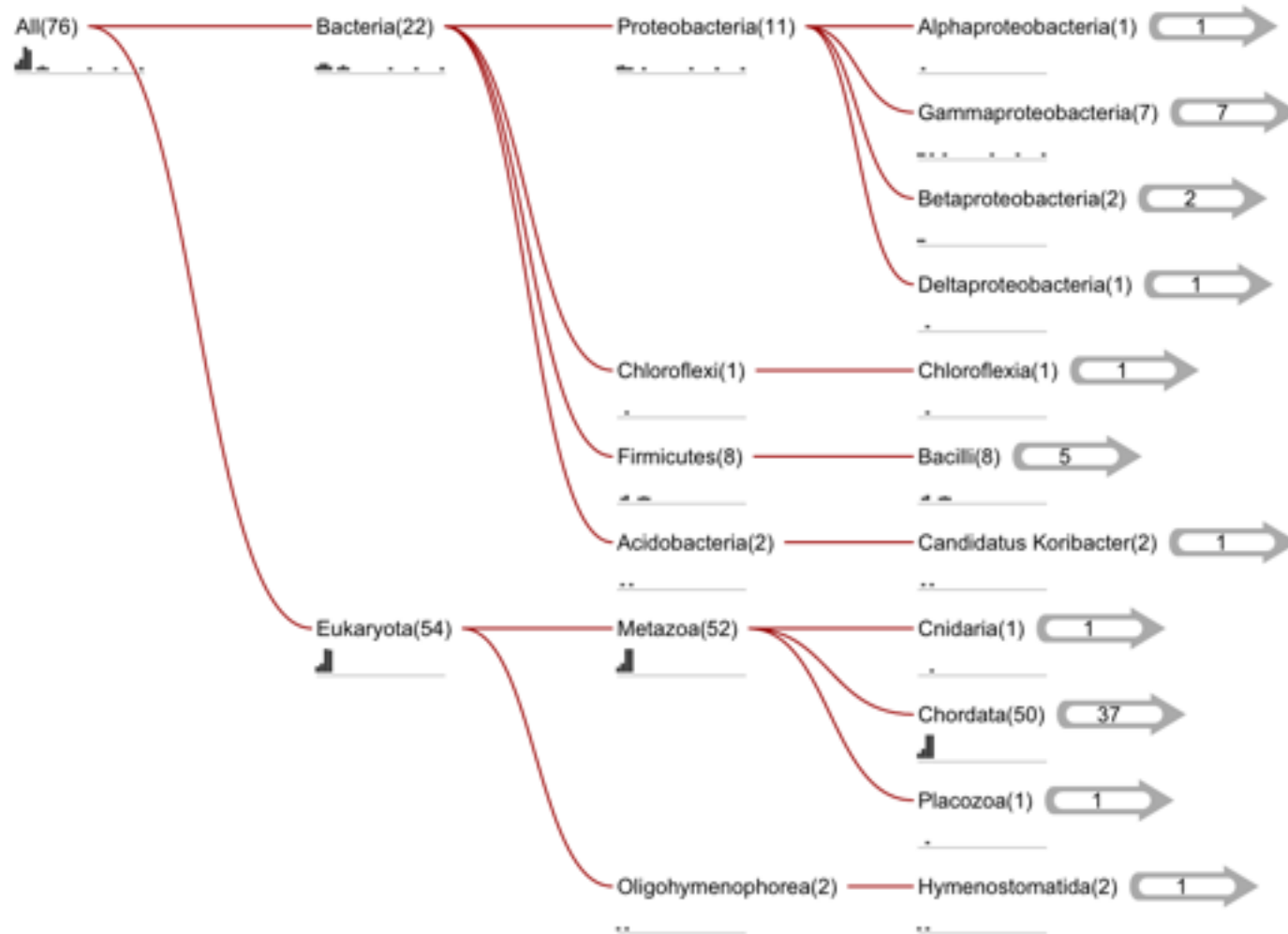
YdjC-like (unpublished data from Lawley Lab)

# Homology and evolution

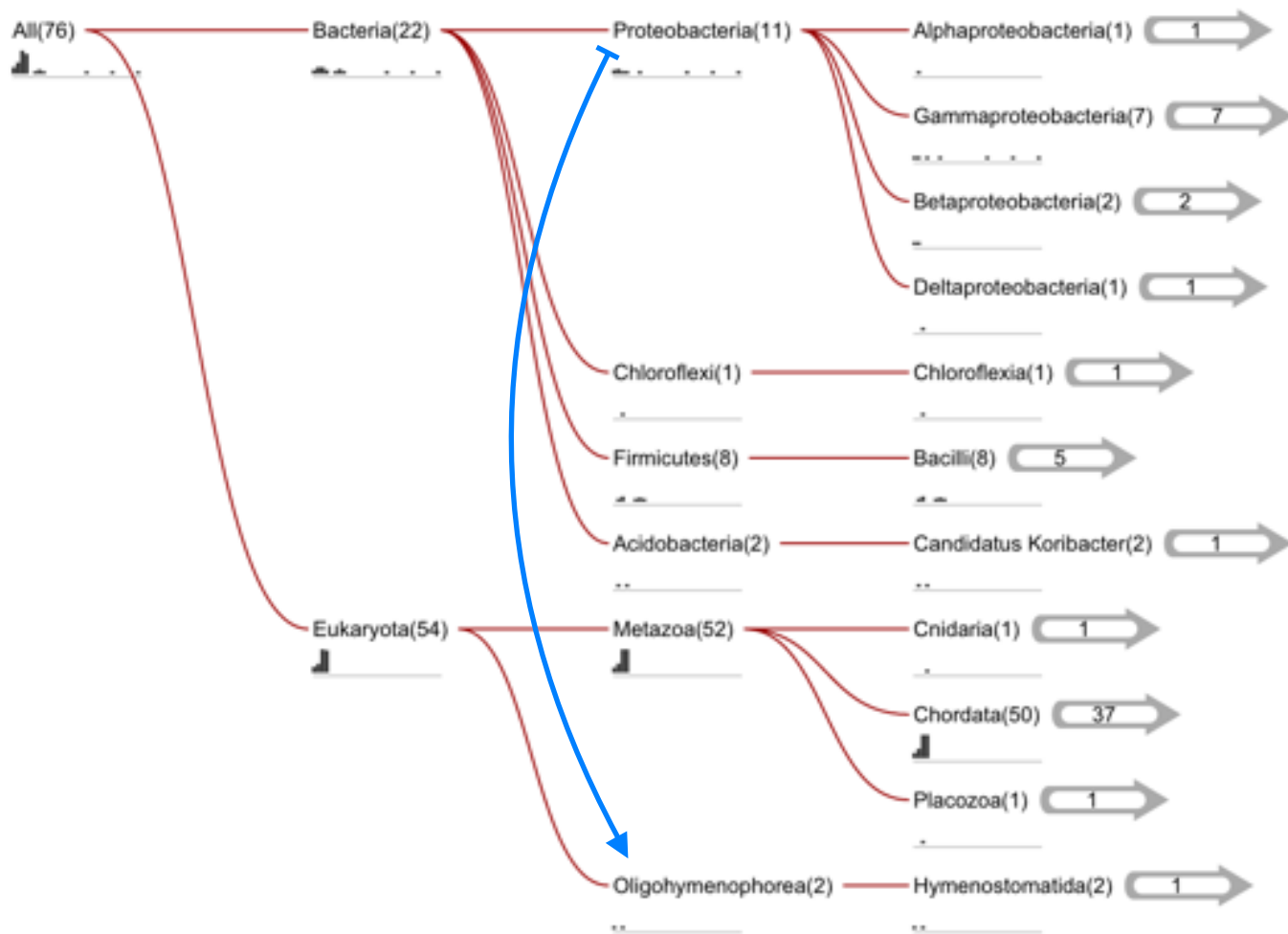
- YdjC gene from *E.coli*
  - Chitooligosaccharide deacetylase



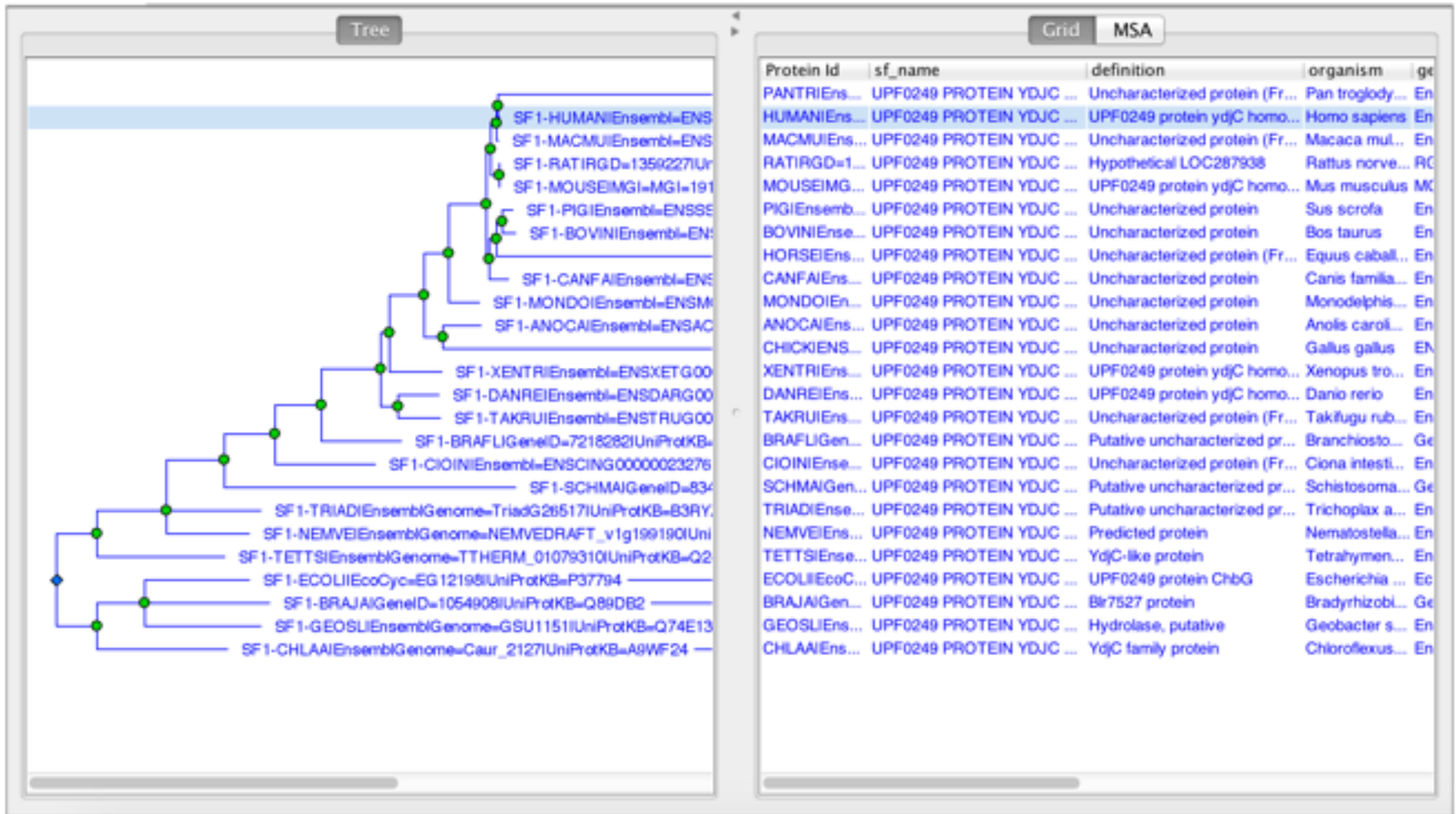
# Homology and evolution



# Homology and evolution



# PANTHER - YdjC



# Acknowledgements

## EMBL-EBI

Alex Mitchell  
Hubert Denise  
Matthew Fraser  
Gift Nuka  
Sebastien Pesseat  
Maxim Scheremetjew  
*Francois Bucchini*  
*Craig McAnulla*  
*Sarah Hunter*

Guy Cochrane  
Rasko Leinonen  
Rajesh Radhakrishnan  
Petra Ten Hoopen

## OeRC

Dawn Field  
Peter Sterk  
  
Susanna Sansone  
Eamonn Maguire  
Alejandra Gonzalez-Beltran  
Philippe Rocca-Serra

## External Collaborators

Sean Eddy  
Eric Nawrocki  
  
Trevor Lawley  
  
Sterghios Moschos  
  
Nils-Peder Willassen

EBI metagenomics - a new resource for the analysis and archiving of metagenomic data

Hunter *et al*, *NAR*, 2014 42:D600-D606

